

# **PyOMP: Writing HPC Code with Python**

Tim Mattson (University of Bristol and Merly.ai ... but mostly retired) tim@timmattson.com

Acknowledgement: The rest of the PyOMP team ....

Giorgis Georgakoudis (LLNL), Todd Anderson (bodo.ai), and Stuart Archibald (anaconda)





# **PyOMP: Writing HPC Code with Python**

Tim Mattson (University of Bristol and Merly.ai ... but mostly retired) tim@timmattson.com

Acknowledgement: The rest of the PyOMP team ....

Giorgis Georgakoudis (LLNL), Todd Anderson (bodo.ai), and Stuart Archibald (anaconda)



# I am best known for my work on OpenMP

I am one of the last members of the original OpenMP team (1996) still active with the language.



<sup>\*</sup> The name "OpenMP" is the property of the OpenMP Architecture Review Board.

# The Growth of Complexity in OpenMP

Our goal in 1997 ... A simple interface for application programmers



The OpenMP specification is so long and complex that few (if any) humans understand the full document

#### How much are different parallel programming models used?

#### HPCorpus data set for training LLM models for parallel programming

Scan all C, C++ and Fortran codes from github with "last updated" dates between 2012 and mid 2023



	Repos	Size(GB)	Files (#)	Functions (#)
С	144,522	46.23	4,552,736	87,817,591
C++	150,481	26.16	4,735,196	68,233,984
Fortran	3,683	0.68	138,552	359,272

Quantifying OpenMP: Statistical insights into usage and adoption, Tal Kadosh, et al., HPEC'2023, https://arxiv.org/abs/2308.08002

HPCorpus Let us assess usage of HPC programming models ... **OpenMP is number 1**!!!



Aggregate numbers over all repositories from 2013 to 2023

Note: since we did not collect files with .cu or .cuf suffixes, we may have undercounted CUDA usage in HPCorpus.

# What are people actually using from OpenMP

With the HPCorpus\* dataset, we finally have hard-data to analyze what "should" be in the common core.

This data was constructed by summing up counts for different directives and clauses across time from 2013 to the middle of 2023.

HPCorpus ... a data set created by scraping "all" HPC codes from github written in C, C++ and Fortran.

Quantifying OpenMP: Statistical Insights into Usage and Adoption, Tal Kadosh, Niranjan Hasabnis, Tim Mattson, Yuval Pinter, and Gal Oren, IEEE HPEC 2023



### **OpenMP CPU Patterns:** Map work onto a team of threads



# **OpenMP GPU Patterns:** Offload work from the CPU to the GPU



#pragma omp target teams loop collapse(2)
for (i=0; i<N; i++)
 for(j=0; j<N; j++)
 for(k=0; k<N; k++)
 C[i][j] += A[B[i][k] \* B[k][j];</pre>

Parallel loops define an index-space. Individual loop iterations define units of work (work-items)



work-items and data mapped onto index space Work-items organized into work-groups



Work-groups enqueued for scheduling.



Work-items from a workgroup execute toether.

# I've been active in HPC since 1980. That makes me an official HPC old-timer

This is my first parallel computer which I used as a Post Doc at Caltech in 1985. We wrote our code using Fortran'77 and a message passing library.

#### The **Caltech Cosmic Cube** developed by Charles Seitz and Geoffrey Fox in1981

- 64 Intel 8086/8087 processors
- 128kB of memory per processor
- 6-dimensional hypercube network



# What HPC old-timers think of Python?

(from the paper, There's plenty of room at the top. Leiserson et. al. Science vol. 368, June 2020).

They used matrix multiplication to explore the connection between software and performance for I in range(4096): for j in range(4096): for k in range (4096): 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	С	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

Amazon AWS c4.8xlarge spot instance, Intel® Xeon® E5-2666 v3 CPU, 2.9 Ghz, 18 core, 60 GB RAM

# What HPC old-timers think of Python

(from the paper, There's plenty of room at the top. Leiserson et. al. Science vol. 368, June 2020).

They used matrix multiplication to explore the connection between software and performance for I in range(4096): for j in range(4096): for k in range (4096): 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	10.0	0.01
3	С	542.67	0.253	Python pe	rformance is a joke.	0.03
4	Parallel loops	69.80	1.969	No seriou	s HPC programmer	0.24
5	Parallel divide and conquer	3.80	36.180	would <b>E</b>	EVER use Python	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

Amazon AWS c4.8xlarge spot instance, Intel® Xeon® E5-2666 v3 CPU, 2.9 Ghz, 18 core, 60 GB RAM

# Why is Python so slow?

• Python is interpreted ... dynamically 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) to allow only one thread at a time can make forward progress.

### Primary Language used in first year, Computer Science Courses



The Reid List tracks a large sample of North American Universities and the languages they use in teaching.

The Reid List was started by Richard Reid in the 1990s. He has retired but others are carrying on the tradition. The above data comes from Trends Of Commonly Used Programming Languages in CS1 And CS2 Learning, Robert M. Siegfried, Katherine G. Herbert-Berger, Kees Leune, Jason P. Siegfried, The 16th International Conference on Computer Science & Education (ICCSE 2021) August 18-20, 2021.

#### **Python is number One!** Popularity of Programming Languages (PyPI)

**Top 5 Languages** 



Programmers have spoken ... Python rules. Old-timers (like me) need to stop being such arrogant snobs and help make Python a first class HPC language

# ... So perhaps best way to bring parallel computing to the masses would be to combine OpenMP and Python?



Multithreaded parallel Python through OpenMP support in Numba Todd Anderson, Timothy G. Mattson, SciPy 2021. http://conference.scipy.org/proceedings/scipy2021/tim\_mattson.html

PyOMP: Multithreaded Parallel Programming in Python Timothy G. Mattson, Todd A. Anderson, Giorgis Georgakoudis, Computing in Science and Engineering, IEEE, November/December 2021

PyOMP: Programming GPUs with OpenMP and Python

Giorgis Georgakouis, Todd A. Anderson, Stuart Archibald, Bronis de Supinski, and Timothy G. Mattson. High Performance Python for Science at Scale workshop at SC24, 2024

# **Design Requirements**



To be an effective parallel computing solution for the Python community, PyOMP must satisfy three requirements

- 1. It must be Pythonic.
  - It must match the way Python programmers working on scientific computing problems use Python.
  - It cannot change Python syntax
- 2. It must deliver performance that is on par with what you'd get from C and OpenMP
- 3. It must be ubiquitous ... available and easy to install on any commonly used platform

# Pythonic OpenMP in three-part harmony

 Incorporated into the numba JIT compiler. The code is JIT'ed into LLVM and therefore avoids the Global Interpreter Lock (GIL) and supports parallel computing with multiple threads.

• Numpy is the standard module used in scientific computing with Python. Hence, PyOMP is optimized to with numpy arrays.

• OpenMP managed through a context manager (that is, a with statement).



**OpenMP** 



# PyOMP by example ...

We will understand PyOMP by considering the three fundamental design patterns of OpenMP (Loop parallelism, SPMD, and divide and conquer) applied to the following problem



# **Loop Parallelism code**

from numba import njit from numba.openmp import openmp\_context as openmp

@njit

def piFunc(NumSteps): step = 1.0/NumSteps sum = 0.0 OpenMP constructs managed through an *openmp* context manager.

Pass the OpenMP directive into the

OpenMP context manager as a string

#### with openmp ("parallel for private(x) reduction(+:sum)"):

for i in range(NumSteps):

x = (i+0.5)\*step sum += 4.0/(1.0 + x\*x)

pi = step\*sum return pi

pi = piFunc(10000000)

Python's implicit data management mapped onto OpenMP. Default rules:

- Variables referenced outside the OpenMP construct are shared
- Variables that only appear inside a construct are private
- Python for technical applications typically based on Numpy arrays, so PyOMP focusses on numpy arrays as well.

OpenMP data environment clauses are supported in PyOMP

# 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
                                         Deal out loop iterations as if a deck of cards (a cyclic distribution)
  return step*np.sum(partialSums)
                                         ... each threads starts with the Iteration = ID, incremented by the
                                         number of threads, until the whole "deck" is dealt out.
pi = piFunc(10000000)
```

# **Divide and Conquer**

 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 only at the leaves
  - Do work as you recombine

# **Divide and conquer (with explicit tasks)**





# Loop Parallelism code naturally maps onto the GPU

from numba import njit from numba.openmp import openmp\_context as openmp

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

with openmp ("target teams loop private(x) reduction(+:sum)"):
 for i in range(NumSteps):
 x = (i+0.5)\*step

 $sum += 4.0/(1.0 + x^*x)$ 

pi = step\*sum return pi

pi = piFunc(10000000)

OpenMP constructs managed through the *with* context manager.

Map the loop onto a 1D index space ... the loop body defines the kernel function

# **Design Requirements**



To be an effective parallel computing solution for the Python community, PyOMP must satisfy three requirements

- It must be Pythonic.
  - It must match the way Python programmers working on scientific computing problems use Python.
  - It cannot change Python syntax
  - 2. It must deliver performance that is on par with what you'd get from C and OpenMP
  - 3. It must be ubiquitous ... available and easy to install on any commonly used platform

We get a "check minus" on programmability since many Python programmers avoid loops and express algorithms solely through numpy array expressions.

PyOMP needs to support the OpenMP workshare construct and implement it with fusion and elision of temporary arrays (something we know how to do based on work on Parallel Accelerator)

# What about performance? Multiple threads running slow, dynamically compiled python code is still slow

# But with statically compiled JIT'ed code ... PyOMP programs are fast.

# Numerical Integration results in seconds ... lower is better

	PyOMP			C/OpenMP		
Threads	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<sup>8</sup> steps

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

For the C programs we used Intel<sup>®</sup> icc compiler version 19.1.3.304 as icc -qnextgen -O3 –fiopenmp

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



# **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) tlnit = 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-8if 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)



Intel<sup>®</sup> Xeon<sup>®</sup> E5-2699 v3 CPU, 18 cores, 2.30 GHz, threads mapped to a single CPU, one thread/per core, first 16 physical cores. Intel<sup>®</sup> icc compiler ver 19.1.3.304 (icc –std=c11 –pthread –O3 xHOST –qopenmp)

# **5-point stencil: Heat diffusion problem**

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

```
# 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
```

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



# **5-point stencil: solve kernel**

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

# **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

22.6 secs

Nvidia V100:

## **Data Movement dominates...**

Solution: parallel stencil (heat)	25,000x25,00 grid for 1	<ul> <li>25,000x25,00 grid for 10 time steps</li> <li>Xeon Platinum 8480+: 67.6 secs</li> </ul>		
<b>@njit</b>	Nvidia V100:	22.6 secs		
<pre>def solve(n, alpha, dx, dt, u, u_tmp):</pre>				
"""Compute the next timestep, given the current	nt timestep"""			
# Finite difference constant multiplier				
r = alpha * dt / (dx ** 2)				
r2 = 1 - 4 * r				
<pre>with openmp ("target loop collapse(2) map(tof: # Loop over the nxn grid</pre>	rom: u, u_tmp)"):	There can be	many time steps	
for i in range(n):				
for j in range(n): For eac			ep, (2*N <sup>2</sup> )*sizeof(TYPE)	
u_tmp[j, i] = (r2 * u[j, i] + bytes move b			etween the host and	
(u[j, i+1] if i < 1	n-1 else 0.0) +	the device		
(u[j, i-1] if i > 0	0 else 0.0) +	L		
(u[j+1, i] if j < 1	n-1 else 0.0) +			
(u[j-1, i] if j > 0	0 else 0.0))			

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

## **Solution:** Explicitly manage the device data environment

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

# **PyOMP HECBench GPU results**

Program	Description		Input	
adam floydwarshall haccmk hotspot3D lavaMD miniBUDE	Adaptive moment estimation Floyd-Warshall path finding a HACC cosmology code micro Thermal modeling for 3D into Molecular dynamics Ligand-protein docking	stochastic optimization for machine learnin algorithm o-kernel egrated circuits stencil computation	10000 200 100 1024 100 16 1000 512 8 5000 power_512x8 temp_512x8 -boxes1d 30 -deck data/bm1 -wgsize 256 -iterations 100	
	adam	floy dwarshall	hotspot3D	
Lime (ns) 0	114.1 110.8 omp pyomp	(m) 20 29.9 26.0	(m) 50 0 70.6 54.8	
	haccmk	lavaMD	minibude	
Time (ms)	2.3 2.9	(su) 100 0 154.0 155.8	() 20 emiliar 20 0 31.3 20.7	

Details in an IWOMP'25 paper submission

AMD EPYC 7763 CPU with an NVIDIA A100 GPU with 80 GB or memory. Python 3.9.18, Numba 0.57, Ilvm-lite 0.40, CUDA 12.2 with driver version 525.105.17

# **Design Requirements**



To be an effective parallel computing solution for the Python community, PyOMP must satisfy three requirements

- 1. It must be Pythonic.
  - It must match the way Python programmers working on scientific computing problems use Python.
  - It cannot change Python syntax
- 2. It must deliver performance that is on par with what you'd get from C and OpenMP
  - 3. It must be ubiquitous ... available and easy to install on any commonly used platform
## How did we implement PyOMP?

# We build on established tools following standard practice in the Python Community

## **PyOMP implementation: CPU**



# **PyOMP** implementation: **CPU** + **GPU**



## **PyOMP:** a Numba extension for upgradeability and maintainability

- Depends on Numba as a compiler toolkit
  - Similar to numba-cuda, numba-hip
- Uses Numba's LLVM dependencies
  - Ilvmlite: provides python bindings for the LLVM API (Currently supports LLVM 14.x We may need to patch PyOMP when Numba moves to LLVM 18/19)
- Tested with Numba 0.57.x, 0.58.x
  - Architectures: linux-64 (x86\_64), osx-arm64 (mac), linux-arm64, linux-ppc64le

PyOMP piggybacks on the off-the-shelf Numba ecosystem.

We don't need to do any extra work to adapt as new versions of Numba are released

# **PyOMP** is easy to install and use

Conda one-line installation

conda install -c python-for-hpc -c conda-forge pyomp

- PyPi package is underway pip install pyomp
- Fast ways to try
  - Binder: <u>https://mybinder.org/v2/gh/Python-for-HPC/binder/HEAD</u>
  - Docker: docker pull ghcr.io/python-for-hpc/pyomp:latest

Open Source code on github:

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

# **Design Requirements**



To be an effective parallel computing solution for the Python community, PyOMP must satisfy three requirements

- 1. It must be Pythonic.
  - It must match the way Python programmers working on scientific computing problems use Python.
  - It cannot change Python syntax
- 2. It must deliver performance that is on par with what you'd get from C and OpenMP
- 3. It must be ubiquitous ... available and easy to install on any commonly used platform

We get a "check minus" on ubiquity since we only work with Nvidia GPUs.

Nvidia works closely with the Numba community so it was straightforward for us to support their GPUs. There is no reason we can't support AMD GPUs, but it will take a small bit or work to make it happen.

In anticipation of future work on AMD GPUs we are refactoring the PyOMP software to make adding other GPUs much easier. Stay tuned

# **Related work: Other OpenMP API bindings**

- Pythran
  - Transpiles python to C++
  - OpenMP using # comments
- PyKokkos
  - Transpiles python to C++
  - OpenMP through Kokkos abstractions
  - Limited support: parallel\_for, parallel\_reduce, parallel\_scan
- OMP4Py
  - Pure python implementation (Threads with GIL disabled)
  - No GPU support (OpenMP version 3.0)
  - Similar interface to PyOMP
  - Slow

Why PyOMP is so important ... and why all of you should start using it

#### In the early days of parallel computing, we were obsessed with finding the "right" parallel programming environment

ABCPL	CORRELATE	GLU	Mentat	Parafrase2	nC
ACE	CPS	GUARD	Legion	Paralation	PC++ SCUEDIJI E
ACT++	CRL	HAsL.	Meta Chaos	Parallel-C++	SCHEDULE
Active messages	CSP	Haskell	Midway	Parallaxis	DOFT
Adl	Cthreads	HPC++	Millipede	ParC	PUEI
Adsmith	CUMULVS	JAVAR.	CparPar	ParLib++	SDDA.
ADDAP	DAGGER	HORUS	Mirage	ParLin	SHMEM SIMDLE
AFAPI	DAPPLE	HPC	MpC	Parmacs	SIMPLE
ALWAN	Data Parallel C	HPF	MOSIX	Parti	Sina
AM	DC++	IMPACT	Modula-P	pC	SISAL.
AMDC	DCE++	ISIS.	Modula-2*	pC++	distributed smalltalk
AppLeS	DDD	JAVAR	Multipol	PCN	SMI.
Amoeba	DICE.	JADE	MPI	PCP:	SUNIC
ARTS	DIPC	Java RMI	MPC++	PH	Split-C.
Athapascan-0b	DOLIB	javaPG	Munin	PEACE	SR
Aurora	DOME	JavaSpace	Nano-Threads	PCU	Sthreads
Automap	DOSMOS.	JIDL	NESL	PET	Strand.
bb_threads	DRL	Joyce	NetClasses++	PETSc	SUIF.
Blaze	DSM-Threads	Khoros	Nexus	PENNY	Synergy
BSP	Ease .	Karma	Nimrod	Phosphorus	Telegrphos
BlockComm	ECO	KOAN/Fortran-S	NOW	POET.	SuperPascal
C*.	Eiffel	LAM	Objective Linda	Polaris	TCGMSG.
"C* in C	Eilean	Lilac	Occam	POOMA	Threads.n++.
C**	Emerald	Linda	Omega	POOL-T	I readMarks
CarlOS	EPL	JADA	OpenMP	PRESTO	TRAPPER
Cashmere	Excalibur	WWWinda	Orca	P-RIO	uC++
C4	Express	ISETL-Linda	OOF90	Prospero	UNITY
CC++	Falcon	ParLin	P++	Proteus	UC
Chu	Filaments	Eilean	P3L	QPC++	V
Charlotte	FM	P4-Linda	p4-Linda	PVM	V <sub>1</sub> C*
Charm	FLASH	Glenda	Pablo	PSI	Visifold V-NUS
Charm++	The FORCE	POSYBL	PADE	PSDM	VPE
Cid	Fork	Objective-Linda	PADRE	Quake	Win32 threads
Cilk	Fortran-M	LiPS	Panda	Quark	WinPar
CM-Fortran	FX	Locust	Papers	Quick Threads	WWWinda
Converse	GA	Lparx	AFAPI.	Sage++	XENOOPS
Code	GAMMA	Lucid	Para++	SCANDAL	XPC 1
COOL	Glenda	Maisie	Paradigm	SAM	Zounds
		Manifold	-		ZPL

Parallel program environments in the 90's

## Language obsessions: More isn't always better

- The Draeger Grocery Store experiment and consumer choice:
  - Two Jam-displays with coupons for a discount on purchase.
    - 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?



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.

#### In the early days of parallel computing, we were obsessed with finding the "right" parallel programming environment

ABCPL ACE ACT++ Active messages Adl Adsmith ADDAP AFAPI ALWAN AM AMDC AppLeS Amoeba ARTS Athapascan-0b Aurora Automap	CORRELATE CPS CRL CSP Cthreads CUMULVS DAGGER DAPPLE Data Parallel C DC++ DCE++ DDD DICE. DIPC DOLIB DOME DOME DOSMOS.	GLU GUARD HAsL. Haskell HPC++ JAVAR. HORUS HPC HPF IMPACT ISIS. JAVAR JADE Java RMI javaPG JavaSpace JIDL	Mentat Legion Meta Chaos Midway Millipede CparPar Mirage MpC MOSIX Modula-P Modula-2* Multipol MPI MPC++ Munin Nano-Threads NESL	Parafrase2 Paralation Parallel-C++ Parallaxis ParC ParLib++ ParLin Parmacs Parti pC pC++ PCN PCP: PH PEACE PCU PET	pC++ SCHEDULE SciTL POET SDDA. SHMEM SIMPLE Sina SISAL. distributed smalltalk SMI. SONiC Split-C. SR Sthreads Strand.
bb threads	DRL	Joyce	NetClasses++	PETSc	SUIF.
Blaze	DSM-Threads	Khoros	Nexus	PENNY	Synergy
BSP	Ease.	Karma	Nimrod	Phosphorus	l elegrphos
BlockComm C*. "C* in C C** CarlOS	With Choice of with all these	overload in min different option	d what did y ns for parallel p	we accomplish programming?	CGMSG. hreads.h++. readMarks RAPPER
CarlOS —	ErL	JADA	Openwir	PRESIO D DIO	uC++
Casilinere C4	Excalibul	ISETL -L inda	ODE90	Prospero	UNITY
CC++	Falcon	ParLin	P++	Proteus	UC
Chu	Filaments	Eilean	P3L	OPC++	V
Charlotte	FM	P4-Linda	p4-Linda	PVM	ViC*
Charm	FLASH	Glenda	Pablo	PSI	Visifold V-NUS
Charm++	The FORCE	POSYBL	PADE	PSDM	VPE
Cid	Fork	Objective-Linda	PADRE	Quake	Win32 threads
Cilk	Fortran-M	LiPS	Panda	Quark	WinPar
CM-Fortran	FX	Locust	Papers	Quick Threads	W W Winda VENOOPS
Converse	GA	Lparx	AFAPI.	Sage++	XPC
Code	GAMMA	Lucid	Para++	SCANDAL	Zounds
COOL	Glenda	Maisie	Paradigm	SAM	ZPL

Parallel program environments in the 90's

#### In the early days of parallel computing, we were obsessed with finding the "right" parallel programming environment

ABCPL ACE ACT++ Active messages Adl Adsmith ADDAP AFAPI ALWAN AM AMDC AppLeS Amoeba ARTS Athapascan-0b Aurora Automap bb_threads Blaze	CORRELATE CPS CRL CSP Cthreads CUMULVS DAGGER DAPPLE Data Parallel C DC++ DCE++ DDD DICE. Furthermore, chasing the spent m	GLU GUARD HAsL. Haskell HPC++ JAVAR. HORUS HPC HPF IMPACT ISIS. JAVAR JADE engineering is next great pro haking the mode	Mentat Legion Meta Chaos Midway Millipede CparPar Mirage MpC MOSIX Modula-P Modula-2* Multipol MPI Stac gramming mod els we have ad	Parafrase2 Paralation Parallel-C++ Parallaxis ParC ParLib++ ParLin Parmacs Parti pC pC++ PCN PCP: DCP: DCP: DCP: DCP: DCP: DCP: DCP:	pC++ SCHEDULE SciTL POET SDDA. SHMEM SIMPLE Sina SISAL. distributed smalltalk SMI. SONiC Snlit-C. Dent F. rgy
BSP	Ease .	Karma	Nimrod NOW	Phosphorus POET	SuperPascal
C*	ECO Fiffel	KUAIN/Fortran-S	NUW Objective Linda	POET. Polaris	TCGMSG.
C . "C* in C	Filean	Lilac	Occam	POOMA	Threads.h++.
C**	Emerald	Linda	Omega	POOL-T	TreadMarks
CarlOS	EPL	JADA	OpenMP	PRESTO	TRAPPER
Cashmere	Excalibur	WWWinda	Orca	P-RIO	uC++
C4	Express	ISETL-Linda	OOF90	Prospero	UNITY
CC++	Falcon	ParLin	P++	Proteus	UC
Chu	Filaments	Eilean	P3L	OPC++	V
Charlotte	FM	P4-Linda	p4-Linda	PVM	ViC*
Charm	FLASH	Glenda	Pablo	PSI	Visifold V-NUS
Charm++	The FORCE	POSYBL	PADE	PSDM	VPE
Cid	Fork	Objective-Linda	PADRE	Quake	Win32 threads
Cilk	Fortran-M	LiPS	Panda	Quark	WinPar
CM-Fortran	FX	Locust	Papers	Quick Threads	WWWinda
Converse	GA	Lparx	AFAPI.	Sage++	XENOOPS
Code	GAMMA	Lucid	Para++	SCANDAL	XPC
COOL					

Parallel program environments in the 90's

# The end of the crisis

 In the early 90's, the HPC community was fed up with message passing chaos. Driven largely by application developers, we created MPI (version 1.0 released in 1994).



 In the late 90's, the HPC community working in the Accelerated Strategic Computing Initiative (ASCI) used their influence over which HPC systems were purchased to "force" vendor's hands to support a standard for programming shared memory systems. The result was OpenMP (version 1.0 released in 1997).



Portable parallel programming is important for the people who create HPC applications. It took their direct involvement and dedication to create open standards and end parallel programming chaos.

#### The major parallel Programming systems in 2024 ... well at least we have our act together in two cases. 🛞

- In HPC, 3 programming environments dominate ... covering the major classes of hardware.
  - MPI: distributed memory systems ... though it works nicely on shared memory computers.

- **OpenMP**: Shared memory systems ... more recently, GPGPU too.

 CUDA, OpenCL, Sycl, OpenACC, OpenMP ...: GPU programming (use CUDA if you don't mind locking yourself to a single vendor ... it is a really nice programming model)

#### Parallel programming with Python is terribly fragmented

dispy Delegate forkmap forkfun Jobibppmap POSH pp pprocess processing **PyCSP PyMP** Ray remoteD torcp VecPv batchlib Celery Charm4py **PyCUDA** Ramba

Dask Deap disco dispy DistributedPYthon exec\_proxy execnet iPython job stream jug mpi4py NetWorkSpaces PaPy papyrus **PyCOMPSs** PyLinda pyMPI pypar multiprocessing **PyOpenCL** 

pyPastSet pypym pynpvm Pyro Ray Rthread ScientificPython.BSP Scientific.DistrubedComputing.MasterSlave Scientific.MPI SCOOP seppo PySpark Python programmers are locked into the same dystopic world of Star-P HPC in the 90's. superrpy torcpy History suggests that this won't StarCluster get better until the python dpctl applications community arkouda demands (and dedicates PyOMP themselves) to a minimal set of dpnp open, standard solutions

# Conclusion

- Python is the language of choice for most programmers ... so let's stop telling them to learn C/C++ or Fortran to do HPC
- PyOMP lets you write OpenMP code in Python. Try it, you'll like it.
- But we need your help ...
  - We need a user base. Please use it and tell us about your successes and failures.
  - Help drive convergence around a minimal number of open, portable parallel programming environments in Python. We all win if this happens.
  - We need more people to join the PyOMP team and help us grow the technology. For example, I want the OpenMP workshare construct with fusion and array elision. I need someone to work with us to make that happen.



My Greenlandic skin-on-frame kayak in the middle of Budd Inlet during a negative tide 52

# The OpenMP Common Core

THE OPENMP COMMON CORE

#### THE OPENMP COMMON CORE

Making OpenMP Simple Again

Timothy G. Mattson, Yun (Helen) He, and Alice E. Koniges For many years now, we've been teaching the subset of OpenMP that is most commonly used. We call this the **OpenMP Common Core** 

We even wrote a book about it.

The list of items in the common core were determined by experience/anecdote ... we didn't have hard data to drive the analysis.

The OpenMP Common Core #pragma omp parallel void omp set thread num() int omp get thread num() int omp get num threads() double omp get wtime() setenv OMP NUM THREADS N #pragma omp barrier #pragma omp critical #pragma omp for #pragma omp parallel for reduction(op:list) schedule (static [,chunk]) schedule(dynamic [,chunk]) shared(list), private(list), firstprivate(list) default(none) nowait #pragma omp single #pragma omp task #pragma omp taskwait

# To learn more about GPU programming with OpenMP

The latest book on OpenMP ...

A book about how to use OpenMP to program a GPU (focusses on C and C++ ... not Python) SCIENTIFIC AND ENGINEERING COMPUTATION

SERIES

# PROGRAMMING YOUR GPU WITH OPENMP

Performance Portability for GPUs

Tom Deakin and Timothy G. Mattson

#### Extra content ...

#### A tutorial introduction to PyOMP (for programmers new to OpenMP)

# **OpenMP**<sup>\*</sup> **Overview**



\* The name "OpenMP" is the property of the OpenMP Architecture Review Board.

• 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")
```

• A parallel multithreaded "hello world" program with PyOMP



- 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

When I run this program, here is the output.

<pre>from numba.openmp import openmp_context as openm @njit def hello():     with openmp("parallel"):         print("hello")         print("world")</pre>
<pre>@njit def hello():     with openmp("parallel"):         print("hello")         print("world")</pre>
print("world")
•

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

hello

world

DONE

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



PE = Processing Element

Figure from "An Introduction to Concurrency in Programming Languages" by J. Sottile, Timothy G. Mattson, and Craig E Rasmussen, 2010

When I run this program, here is the output.

from numba i	mport niit
from numba.o	penmp import openmp_context as openmp
@njit	
<pre>def hello():</pre>	
with ope	<pre>nmp("parallel"): nt("bello")</pre>
pri	nt("world")
Г <sup>.</sup> . –	
hello()	
print("DONE"	)

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.

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



- 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.

- 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
```

- 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.



- 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.



- 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.



# **Loop Parallelism code**

from numba import njit

from numba.openmp import openmp\_context as openmp

@njit

def piFunc(NumSteps): step = 1.0/NumSteps pisum = 0.0 OpenMP managed through the with context manager.

Numba Just In Time (JIT) compiler compiles the Python code into LLVM thereby bypassing the GIL. Compiled code cached for later use.

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(10000000)

Pass the OpenMP directive into the OpenMP context manager as a string

- 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(10000000)
```
### Numerical Integration results in seconds ... lower is better

	PyOMP C			
Threads	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	

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

For the C programs we used Intel<sup>®</sup> 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



## 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()
                                                         omp_get_num_threads(): get N=number of threads
    with openmp("single"):
                                                     ٠
                                                         omp_get_thread_num(): thread rank = 0...(N-1)
       numThrds = omp get num threads()
                                                     ٠
                                                         single: One thread does the work, others wait
    localSum = 0.0
                                                     •
                                                         private(x): each threads gets its own x
    for i in range(threadID, NumSteps, numThrds):
                                                         shared(x): all threads see the same x
                                                     ٠
      x = (i+0.5)^*step
       localSum = localSum + 4.0/(1.0 + x^*x)
                                                     Deal out loop iterations as if a deck of cards (a cyclic distribution)
    partialSums[threadID] = localSum
                                                     ... each threads starts with the Iteration = ID, incremented by the
  return step*np.sum(partialSums)
                                                     number of threads, until the whole "deck" is dealt out.
```

pi = piFunc(10000000)

## 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.



### Numerical Integration results in seconds ... lower is better

	PyOMP			С		
Threads	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<sup>8</sup> steps

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

For the C programs we used Intel<sup>®</sup> icc compiler version 19.1.3.304 as icc -qnextgen -O3 –fiopenmp

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 <u>without</u> 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.



### **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)**



### Numerical Integration results in seconds ... lower is better

Threads	PyOMP			С		
	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<sup>8</sup> steps

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

For the C programs we used Intel<sup>®</sup> icc compiler version 19.1.3.304 as icc -qnextgen -O3 –fiopenmp

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**

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

### **PyOMP subset of OpenMP for CPU programming**

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

### 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]



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	С	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

Amazon AWS c4.8xlarge spot instance, Intel® Xeon® E5-2666 v3 CPU, 2.9 Ghz, 18 core, 60 GB RAM

### **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) tlnit = 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-8if 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) 89

### DGEMM PyOMP vs C-OpenMP

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



Intel<sup>®</sup> Xeon<sup>®</sup> E5-2699 v3 CPU, 18 cores, 2.30 GHz, threads mapped to a single CPU, one thread/per core, first 16 physical cores. Intel<sup>®</sup> icc compiler ver 19.1.3.304 (icc –std=c11 –pthread –O3 xHOST –qopenmp)

### And we can use PyOMP for GPU programming

### The "BIG idea" Behind GPU programming



#### **Data Parallel vadd with CUDA**



### 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.



 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.

### **SIMT:** One instruction stream maps onto many SIMD lanes

SIMT model: Individual scalar instruction streams are grouped together for SIMD execution on hardware



### A Generic GPU (following Hennessey and Patterson)



### A Generic GPU (following Hennessey and Patterson)



### **A Generic Host/Device Platform Model**



- One *Host* and one or more *Devices*
  - Each Device is composed of one or more Compute Units
  - Each Compute Unit is divided into one or more *Processing Elements*
- Memory divided into *host memory* and *device memory*

PE: processing element. The finest-grained processing element inside a GPU. Also known as a SiMD-lane or CUDA-core.

### **Executing a program on CPUs and GPUs**





One work-group per compute-unit executing

### **Executing a program on CPUs and GPUs**



One work-group per compute-unit executing

SIMD Lanes

SIMD Lanes

### **CPU/GPU execution models**



For a CPU, the threads are all active and able to make forward progress.

For a GPU, any given work-group might be in the queue waiting to execute. 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



Only the arrays are moved back to the host after the target region completes

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

@njit def main(): N = 1024A = 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



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>>>

OpenMP constructs managed through the *with* context manager.

Map the loop onto a 2D index space ... the with openmp("target teams loop collapse(2) private(j)"): loop body defines the kernel function for i in range(N): for k in range(N): target: map execution from the host onto the device for j in range(N):

C[i][j] += A[i][k] \* B[k][j]

- - 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).

```
@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.

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

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

## **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

22.6 secs

Nvidia V100:

### **Data Movement dominates...**

#### 25,000x25,00 grid for 10 time steps

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

22.6 secs



- 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.



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.