JuliaCon 2019 Parallel Computing Workshop

Welcome!

Parallel Computing workshop

- Download JuliaPro version 1.1.1.1
 - <u>https://juliacomputing.com/products/juliapro</u>
- Install it!
- Download the workshop materials
 - github.com/mbauman/ParallelWorkshop2019
- Open Julia Pro, go to File -> Add Project Folder... and select the downloaded folder

Troubleshooting tips

- Ensure running Julia 1.1.
- Running on JuliaPro works best
- Make sure you're in the "Project folder" specifically for the downloaded folder (not one level higher)
- If running on vanilla Juno + Julia, run
]activate
]rm CredentialsHandler
]instantiate

Goals

- Understand modern parallel architectures
- Write and run a multithreaded algorithm on your own computer
- Learn how to structure your program to avoid race conditions
- Treat your computer like a cluster and run multiprocess code on it

What is happening to our computers?



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2015 by K. Rupp

Historical microprocessor trends



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2015 by K. Rupp

Parallelism is required for best performance

Computers today:

- Have single processor cores that can do more than one thing at once
- Have multiple cores/processors (up to 50 cores on a chip, multiple chips)
- Can network multiple computers together
- May have specialized compute hardware (GPU, TPU)

Fastest supercomputers

Rank	Site	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power (kW)
1	DOE/SC/Oak Ridge National Laboratory United States	Summit - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband IBM	2,414,592	148,600.0	200,794.9	10,096
2	DOE/NNSA/LLNL United States	Sierra - IBM Power System S922LC, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband IBM / NVIDIA / Mellanox	1,572,480	94,640.0	125,712.0	7,438
3	National Supercomputing Center in Wuxi China	Sunway TaihuLight - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway NRCPC	10,649,600	93,014.6	125,435.9	15,371

Fastest supercomputers

Rank	Site	System	1608 nodos		
1	DOE/SC/Oak Ridge National Laboratory United States	Summit - IBM Power Sys AC922, IBM POWER9 220 NVIDIA Volta GV100, Dual Mellanox EDR Infiniband IBM	2 processors per node 22 cores per processor 6 GPUs per node		
2	DOE/NNSA/LLNL United States	Sierra - IBM Power Syste S922LC, IBM POWER9 22 NVIDIA Volta GV100, Dual Mellanox EDR Infiniband IBM / NVIDIA / Mellanox	-rail		
3	National Supercomputing Center in Wuxi China	Sunway TaihuLight - Sun Sunway SW26010 260C 1. Sunway NRCPC	way MPP, 10,649,600 93,014.6 125,435.9 15,371 45GHz,		

Parallel Challenges

- Chasing peak performance
 Need a fast language for biggest benefit
- We learn to reason about & write programs serially Need to safely express parallelism
- Often dovetails with large datasets
 Need to consider data movement costs
- Vast array of parallel hardware configurations
 Need to pick the strategy that fits your system

Parallel Strategies

- Making the most of one core: high performance Julia and SIMD
- Making the most of one computer: multithreading
- Making the most of multiple computers/clusters: distributed computing
- Enabling external accelerators:
 GPUs and TPUs

		Ð
	» Manual » Performance Tips	C Edit on GitHub
v1.1.1	Performance Tips	
Search docs	In the following sections, we briefly go through a few make your Julia code run as fast as possible.	r techniques that can help
Performance Tips		
Avoid global variables		
Measure performance with @time	Avoid global variables	
and pay attention to memory allocation	A global variable might have its value, and therefore point. This makes it difficult for the compiler to optim	its type, change at any nize code using global
Tools	variables. Variables should be local, or passed as argu	uments to functions,
Avoid containers with abstract type	whenever possible.	
parameters	Any code that is performance critical or being bench	marked should be inside a
Type declarations	function	

Single-instruction, multiple-data (SIMD)

 $egin{aligned} x_1+y_1 &
ightarrow z_1 \ x_2+y_2 &
ightarrow z_2 \ x_3+y_3 &
ightarrow z_3 \ x_4+y_4 &
ightarrow z_4 \end{aligned}$ four instructions



one instruction!

Single-instruction, multiple-data (SIMD)



1• f	or	x	in	xs	
2		#			
3 ^ e	end				

• Happens automatically!

Single-instruction, multiple-data (SIMD)





- Happens automatically (if it's safe)!
 - Requires "straight line" code: no if statements, @inbounds array accesses, no breaks, only call "small" functions, etc.
 - Limited to basic processor instructions (simple arithmetic)

Single-instruction, multiple-data (SIMD)





- **@simd** can help it happen more frequently
 - Allows floating point re-associativity
- Using smaller data types can increase parallelism
 - Up to 512 bits can process 32 Float32s at a time!

Understanding data movement costs

	Actual cost Scaled "huma cost	
One CPU Cycle	0.4 ns	1 s
Level 1 cache access	0.9 ns	2 s
Level 2 cache access	2.8 ns	7 s
Level 3 cache access	28 ns	1 min
Main memory access	~100 ns	4 min
NVMe SSD I/O	~25 µs	17 hours
SSD I/O	50-150 µs	1.5-4 days
Rotational disk I/O	1-10 ms	1-9 months
Internet call: SF to NYC	65 ms	5 years

https://www.prowesscorp.com/computer-latency-at-a-human-scale

Multithreading

Hands on!

Parallel Algorithm Design

Demo

Tasks

Quick demo

Distributed Computation

Towards clusters!

Mental model



Distributed demo!

GPUs and more

What's a CPU look like?



What's a GPU look like?





What's a GPU look like?



- 15 "multiprocessors"
- Up to 6 warps per multiprocessor

What's a GPU look like?



- 15 "multiprocessors"
- Up to 6 warps per multiprocessor
- Each warp is 32 threads



Warp programming





A();B(); X();Y();

Time

GPUs: Using JuliaGPU to make it easy

julia> cu(rand(100))
100-element CuArray{Float32,1}:
0.5351649
0.93891734
0.6757371
0.4985685
0.9094223
0.8695512
0.0144788
0.77135414
0 612888

- Move data to the GPU
- CuArray is specialized to operate on the GPU and favors Float32s
- In particular, you can create your own "GPU kernels" with broadcast fusion
- You can also manually write CUDA code



TPU: the Tensor Programming Unit

Scales to pods (512 TPU cores - 4.3 PF_{16} /s on ResNet50)



Fischer et al. Automatic Full Compilation of Julia Programs and ML Models to Cloud TPUs (<u>arXiv:1810.09868</u>)

Julia on TPUs

- Unchanged Language Semantics on TPUs
 - Control Flow
 - Multiple Dispatch
 - Data Structures
- Share code between CPU/GPU/TPU
- Integration with native distributed computing facilities (WIP)
- Goal: Make TPU pods and traditional HPC systems look the same to the programmer. Seamless retargeting.

Happy parallel computing Happy parallel computing Happy parallel computing Happy parallel computing