

Porting MURaM (Max Planck University of Chicago Radiative MHD) to GPUs Using OpenACC

Rich Loft (Director TDD in CISL, NCAR) Eric Wright (PhD student) & Sunita Chandrasekaran (Assistant Professor) University of Delaware <u>loft@ucar.edu</u>, {efwright, <u>schandra}@udel.edu</u>

Project in collaboration with NCAR, Max Planck for Solar System Research and University of Delaware



March 19, 2019 GTC 2019



Anatomy of the Sun

ELAWARE.







💿 NVIDIA.

GPU EDUCATION

CENTER

2

MURaM (Max Planck University of Chicago Radiative MHD)

- The primary solar model used for simulations of the upper convection zone, photosphere and corona.
- Jointly developed and used by HAO, the Max Planck Institute for Solar System Research (MPS) and the Lockheed Martin Solar and Astrophysics Laboratory (LMSAL).
- MURaM has contributed substantially to our understanding of solar phenomena.
- MURaM also plays a key role in interpreting high resolution solar observations.

The Daniel K. Inouye Solar Telescope (DKIST), a ~\$300M NSF investment, is expected to advance the resolution of ground based observational solar physics by an order of magnitude.



MURaM simulation of solar granulation



EDUCATION

CENTER





High Resolution Simulations of the Solar Photosphere

- Forward modeling of DKIST observables will require simulations with grid spacing of 4 km on a regular basis.
- Requires at least 10-100x increase in computing power compared to current baseline.









From data inspired to data driven simulations of solar eruptions

- Realistic simulations of the coupled solar atmosphere are an important tool to understand and even predict solar eruptions.
- Current models run about ~100x slower than real-time
- Data driven simulations of solar events would allow for analysis and prediction of ongoing solar events
- Future data assimilation applications will require ensemble runs (~10x)



Comprehensive model of entire life cycle of a solar flare (Cheung et al 2018)





Realistic Simulations of the Chromosphere

- The chromosphere is a difficult region to observe and model
- New instruments (DKIST 4m telescope and sunrise balloon observatory) will allow for unprecedented observations
- However, the modeling still must be brought up-to-date to these observational advancements.







Toward Predicting Solar Eruptions

- How do we get MURaM 100x faster? GPU-accelerated RT is the key.
 - RT solver speed up on GPUs (~85x) a CPU core via <u>wavefront algorithms</u> (Chandrasekaran et al.)
 - RT iteration counts are reduced (~2x) on GPU's bigger subdomains.
 - But, we can rewrite RT solver to do required wavelengths (these are embarrassingly parallel).
 - Last two points play to GPU's strength: data parallelism.
 - Estimate 450 GPUs could achieve breakeven (1 simulated hour/hour)
- Thousands of GPUs would be required to do full data assimilation.
 - Expensive but not unthinkable (Summit has 27,600!)
- Requirements play to strengths of GPUs, and trend in their design.







Planned MURaM Development

- Porting of the MURaM code to GPUs using OpenACC (collaboration between HAO, CISL, University of Delaware and MPS)
- Implementation of a more sophisticated chromosphere (HAO, MPS)
- Implementation of boundary conditions that allow data driven simulations of solar events (HAO, CU Boulder, LMSAL)
- New IO modules that allow for data compression during the IO process and runtime visualization.







Time to dive deeper into the computational science side of things! ③

ELAWARE







- Profiling, parallelization, re-profiling
- Optimizing CPU-GPU data transfer/management
- Focusing on the most important loop Radiation Transport (RT)
 - Long term science goals
- Re-designing Radiation Transport Algorithm
- More profiling info to find & address performance challenges





Νυσιλα

GPU EDUCATION CENTER

12



GAUSSIAN 16



Image coardeay ANSIYS

ANSYS FLUENT

We've effectively used **OpenACC** for heterogeneous computing in ANSYS Fluent with impressive performance. We're now applying this work to more of our models and new platforms.

NUMECA FINE/Open





Porting our unstructured C++ CFD solver FINE/Open to GPUs using OpenACC would have been impossible two or three years ago, but OpenACC has developed enough that we're now getting some really good results.



SYNOPSYS

For VASP, OpenACC is the way

forward for GPU acceleration.

cases better than CUDA C. and

OpenACC dramatically decreases

with NVIDIA and PGI as an early

adopter of CUDA Unified Memory.

GPU development and maintenance

efforts. We're excited to collaborate

Performance is similar and in some

VASP

Using OpenACC, we've GPUaccelerated the Synopsys TCAD Sentaurus Device EMW simulator to speed up optical simulations of image sensors. GPUs are key to improving simulation throughput in the design of advanced image sensors.



OpenACC made it practical to develop for GPU-based hardware while retaining a single source for almost all the COSMO physics code

"

MPAS-A



Our team has been evaluating OpenACC as a pathway to performance portability for the Model for Prediction (MPAS) atmospheric model. Using this approach on the MPAS dynamical core, we have achieved performance on a single P100 GPU equivalent to 2.7 dual socketed Intel Xeon nodes on our new Cheyenne supercomputer.



With OpenACC and a compute node based on NVIDIA's Tesla P100 GPU, we achieved more than a 14X speed up over a K Computer node running our earthquake disaster simulation code



Adding OpenACC into MAS has given us the ability to migrate medium-sized simulations from a multi node CPU cluster to a single multi-GPU server. The implementation yielded a portable single-source code for both CPU and GPU runs. Future work will add OpenACC to the remaining model features, enabling GPU-accelerated realistic solar storm modeling.



Image country Oak Ridge National Laboratory





E3SM

The CAAR project provided us with

early access to Summit hardware and

access to PGI compiler experts. Both

of these were critical to our success.

PGI's OpenACC support remains the best available and is competitive with

much more intrusive programming

model approaches.

Due to Amdahi's law, we need to port more parts of our code to the GPU if we're

SANJEEVINI

going to speed it up. But the sheer number of routines poses a challenge. OpenACC directives give us a low-cost approach to getting at least some speedup out of these second-tier routines. In many cases it's completely sufficient because with the current algorithms, GPU performance is bandwidth-bound.





Using OpenACC our scientists were able to achieve the acceleration needed for integrated fusion simulation with a minimum investment of time and effort in learning to program





GPUs.



DpenAGC can prove to be a handy tool for computational engineers and researchers to obtain fast solution of non-linear dynamics CED, we have obtained order of magnitude reduction in computing time by porting several Especially the routines involving search algorithm and matrix solvers have been well-accelerated to improve the overall scalability of the code.



PWscf (Quantum ESPRESSO)

> CUDA Fortran gives us the full performance potential of the CUDA programming model and NVIDIA GPUs. While leveraging the potential of explicit data movement, ISCUF KERNELS directives give us productivity and source code maintainability. It's the best of both worlds.



mage courtesy: NCAR







"





In an academic environment maintenance and speedup of existing codes is a tedious task. OpenACC. provides a great platform for computational scientists to accomplish both tasks without involving a lot of efforts or manpower in speeding up the entire computational task.



"







<

14



Profiling Tools

- Several profiling techniques were used to obtain an initial, high-level view of the code
- Function call map
- Arm MAP for generalized performance metrics and MPI
- NVProf for GPU performance profiling









Experimental Setup

- NVIDIA PSG Cluster (hardware)
 - CPU: Intel Xeon E5-2698 v3 (16-core) and Xeon E5-2690 v2 (10-core)
 - GPU: NVIDIA Tesla V100 (4 GPUs per node)
- Software Used
 - PGI 18.4 (CUDA 9.1 and 9.2) and PGI 18.10
 - Results in this talk use PGI 18.4
 - icc 17.0.1/18.0.1





ΔΙΟΙΛΟ 🗵

GPU EDUCATION CENTER

15



Routine Descriptions

GPU
EDUCATION CENTER

Name	Routine Summary:	Broadwell (v4) core (sec)
TVD Diffusion	Update diffusion scheme - using TVD slope + flux limiting.	7.36812
Magnetohydrodynamics	Calculate right hand side of MHD equations.	6.26662
Radiation Transport	Calculate radiation field and determine heating term (Qtot) required in MHD.	5.55416
Equation of State	Calculate primitive variables from conservative variables. Interpolate the equation of state tables.	2.26398
Time Integration	Performs one time integration.	1.47858
DivB Cleaner	Clean any errors due to non-zero div(B).	0.279718
Boundary Conditions	Update vertical boundary conditions.	0.0855162
Grid Exchange	Grid exchanges (only those in Solver)	0.0667914
Alfven Speed Limiter	Limit Maximum Alfven Velocity	0.0394724
Synchronize timestep	Pick minimum of the radiation, MHD and diffusive timesteps.	4.48E-05

NCAR











🔰 NVIDIA

Profile driven parallelization

- Based on information gathered from profiling, implement simple development cycle:
 - Identify which loops are currently the most impactful
 - Parallelize the loop(s) for GPU execution
 - Verify that our test cases pass with the new change
 - Reprofile/Optimize until happy enough to move







💿 NVIDIA.

GPU Profile using nvprof







CUDA Occupancy Report

FIVERSITY OF

240x160x160 Dataset

NCAR

Kernel Name	Theoretical Occupancy	Achieved Occupancy	Runtime % (GPU)
MHD	25%	24.9%	32.4%
TVD	31%	31.2%	31.6%
CONS	25%	24.9%	6.3%
Source_Tcheck	25%	24.9%	5.2%
	Radiation	Transport	
Driver	100%	10.2%	15.2%
Interpol	56%	59.9%	4.9%
Flux	100%	79%	1.5%





What did we learn so far?

- What is MURaM?
- What is the state of the project?
- What are the challenges identified?
- What problems do we still have to overcome?
 - Optimizing RTS
 - Learning from CUDA Occupancy Report







Toward Predicting Solar Eruptions

- How do we get MURaM 100x faster? GPU-accelerated RT is the key.
 - RT solver speed up on GPUs (~85x) a CPU core via <u>wavefront algorithms</u> (Chandrasekaran et al.)
 - RT iteration counts are reduced (~2x) on GPU's bigger subdomains.
 - But, we can rewrite RT solver to do required wavelengths (these are embarrassingly parallel).
 - Last two points play to GPU's strength: data parallelism.
 - Estimate 450 GPUs could achieve breakeven (1 simulated hour/hour)
- Thousands of GPUs would be required to do full data assimilation.
 - Expensive but not unthinkable (Summit has 27,600!)
- Requirements play to strengths of GPUs, and trend in their design.





Data Dependency Along Rays:

D B

Figure 4.2: The walking order of the Short Characteristics method in a 2D grid for a ray direction pointing into the upper right quadrant. Black circles represent gridpoints on the upwind boundaries, where the intensity values are assumed to be known.

- Data dependency is along a plane for each octant, angle combo.
- Depends on resolution ratio, not known until run-time.
- Number of rays per plane can vary.

Vögler, Alexander, et al. "Simulations of magneto-convection in the solar photosphere-Equations, methods, and results of the MURaM code." Astronomy & Astrophysics 429.1 (2005): 335-351.

Figure 4.1: The intensity at gridpoint F is obtained by solving the transfer equation along the short characteristic \overline{EF} . The intensity at the upwind point E is interpolated from the (already known) intensity values at the surrounding gridpoints, A to D.











Solving RTS Data Dependency

- We can deconstruct the 3D grid into a series of 2D slices
- The direction of the slices is dependent on the X,Y,Z direction of the ray
- Parallelize within the slice, but run the slices themselves serially in predetermined order









💿 NVIDIA.













💿 NVIDIA.













💿 nvidia.

*driver.nvvp	¶ *gpu_i_n2.nvvp ∞					•	
9016.995 ms	9016.9975 ms	9017 ms	9017.0 <mark>4.73 µs</mark>	9017.005 ms	9017.0075 ms	9017.01 ms	
	acc_compute	construct@rt.cc:12	42	acc_comp	ute_construct@rt.cc:1242		
	acc_enque	ue_launch@rt.cc:124	42	acc_en	iqueue_launch@rt.cc:1242	2	
		cuLaunchKernel			cuLaunchKernel		
	RTS::driver_1242_	gpu(double		RTS::driv	er_1242_gpu(doubl		
	RTS::driver_1242_	gpu(double		RTS::driv	er_1242_gpu(doubl		_







What did we learn so far?

- What is MURaM?
- What is the state of the project?
- What are the challenges identified?
- What problems do we still have to overcome?
 Optimizing RTS
- Learning from CUDA Occupancy Report







NVIDIA V100 specification

- 64 warps per multiprocessor (32 threads per warp)
- 65536 registers per multiprocessor
- 96 KB shared memory per multiprocessor
- Occupancy = # Used Warps / # Possible Warps
- How many warps we can use is dependent on how many registers and how much shared memory we use per thread

Thanks to NVIDIA for giving us access to their PSG system for our experiments!!!





CUDA Occupancy Report

- Using the CUDA occupancy calculator
- We can see why our theoretical occupancy for MHD, TVD, etc is so low
- From this graph, we know that threads-per-block is not the problem









CUDA Occupancy Report

- We can also see that our shared memory usage is much lower than it could be
- One idea we have is to start moving some of our data to shared memory to get better usage



CUDA Occupancy Report

- Finally, we see that our problem is register usage per thread
- To get 100% occupancy, we would need to reduce register usage to 32 registers per thread

MURaM: Parallel processing using MPI

Introduction

- Three dimensional uniform cartesian grid.
- Grid is divided into smaller pieces based on MPI ranks and are processed independently.
- Halo information is communicated between the MPI ranks at a regular intervals.
- Large percentage of halo exchange happens in Radiative transfer.

Optimization strategies

- Overlap computation and communication in RTS.
- Perform a GPU to GPU direct communication without involving host.
- Optimize the packing and unpacking of communication data.

💿 nvidia.

MURaM: MPI profiling using Intel Trace Analyzer

NIVERSITYOF

4.84 s	4.92 s	5.00 s	5.08 s		5.16 s		5.24 s	5	.32 s	5.40 s	5	.48 s
PO AdAprApplid	18 S 4 Millellellellellellellellellellellellelle			5.1/2 S nlication Δηροίι	5.20 W4000lication		5.28 S	Annli			4 S Δnnlication	سمامہ مامہ مام
	i <i>u, LE LE LE LE EL LE LE LE LE LE LE</i>			cation (Appli		Appl		Appli			Application	
	i <i>uu ka ka</i>	14		plicition (Appli		Appl	AppliApplication	Appli	AppliAApplication		Application	
	ive te	44 44 44 44 44 44 44 44 44		plica on Appli		Appla	AppliApplication	JilgaA	AppliAApplication		Application	
	i <i>wa<mark>i</mark> ka ka</i>	14	: <i>44, 44, 44</i> : 44 Ap Appl A Ap	plication (Appli	AppliApplication	Appl.	AppliApplication	Appli	AppliAApplication	Applic & Application	Application	
P5 A Ap Appliv AAppl	ive de ee, ee, ee, ee, ee, ee, ee, ee, ee,	14. 14. 14. 14. 14. 14. 14. 14. 14. 14.	144, 44, 44, 44 Ap Appl Al Ap	plicat n ⁴ Appli	AppliApplication	Appl.	AppliApplication	Appli	AppliAApplication	Applic Application	Application	
P6 A Ap Appli/MAAppl	ive de ee, ee, ee, ee, ee, ee, ee, ee, ee,	14. 44, 14. 44. 44 <mark>. 44, 14. 44, 14.</mark> 44,	144, 44, 44, 44 Ap Appl A Ap	plicat n ⁴ Appli	AppliApplication	Appl.	AppliApplication	Appli	AppliAApplication	Applic Application	Application	
P7 A <mark>,</mark> ∕Ap Appli⊮ <mark>M</mark> AAppl	1144	14. 14. 14. 14 <mark>. 14.</mark> 14. 14. 14. 14.	: <i>44<mark>, 44, 44:</mark> 4</i> 4 Ap Appl A <mark>l</mark> Apl	plicat n ⁷ Appli	AppliApplication	Appl	AppliApplication	Appli	AppliAApplication	Applic SApplication	Application	<u>, AA, AA, AA, </u>
P8 A <mark>. Ap appli/M</mark> AAppl	ive. IEE EE, EE, EE, EE, EE, EE, EE, EE, EE,	14. 14. 14. 14. 14. 14. 14. 14. 14.	: <i>44<mark>,</mark> 44, 44</i> ; 44Ар Арр А <mark>:</mark> Арј	plicat n 🦊 Appli	AppliApplication	Арр/	AppliApplication	Арри	AppliAApplication	Applic & Application	Application	
P9 A <mark>l Apappli/M</mark> AAppl	ì <i>dh<mark>a</mark> 44</i> , 44, 44, 44, 44, 44, 44, 44, 44, 44	14, 14, 14, 14, 14, 14, 14, 14, 14, 14,	: <i>44, 44, 44</i> ; 444p App 4; Apj	plicat n <mark>/</mark> Appli	AppliApplication	Арр. <mark>М</mark>	AppliApplication	Арр	AppliAApplication	MApplic Application	Application	
P10A <mark>.,</mark> Ap appli/ <mark>M</mark> AAppl	1 <i>44, 1</i> 44, 44, 44, 44, 44, 44, 44, 44, 44, 44,	14. 14. 14. 14. 14. 14. 14. 14. 14. 14.	: <i>44, 44, 44</i> : 44АрАрр А <mark>:</mark> Арј	plicat n <mark>/</mark> Appli	AppliApplication	App	AppliApplication	Арр	AppliAApplication	MApplic CApplication	Application	AAAAAA
P11A <mark>,</mark> Ap appli/ <mark>M</mark> AAppl	114 - 1 42 44, 44, 44, 44, 44, 44, 44, 44, 44, 4	14. 14. 14. 14. 14. 14. 14. 14. 14. 14.	. 44, 44, 44, 44 4 Ap App A Ap	plicat n 🤼 🗛 pli	AppliApplication	App/ <mark>/</mark>	AppliApplication	Арри	AppliAApplication	MApplic - Application	Application	
P12A <mark>, Ap Applica AAppl</mark>	iva <mark>n</mark> 146, 44, 44, 44, 44, 44, 44, 44, 44, 44,	14. 44, 14. 44. 14, 14, 14, 14, 14, 14,	44, 44, 44, 44 Ap App A Ap	plicat n ⁷ Appli	ApplicatiApplicatio	n 🛛 🖌	ApplicaApplication	App	Applic&AApplication	Mapplic Application	Application	אי אאי אאי
P13A <mark>, Ap Applica A</mark> Appl	A, , , , , , , , , , , , , , , , , , ,	14. 44, 14. 44. 14, 14, 14, 14, 14, 14,	44, 44, 44, 44 Ap App A Ap	plicat n ⁽ Appli	ApplicatApplication	n 시	ApplicatApplication		ApplicAApplication	Applic Application	Application	AA AA AA
P14A Ap Applica AAppl	14 FE,	14, 14, 14, 14, 14, 14, 14, 14, 14, 14,	44, 44, 44, 44 Ap App A Ap	plicat n ⁽ Appli	ApplicaApplication	A <mark>¢.</mark> ¢	ApplicatApplication	A	Applic:A;Application	Applic •• Application	Application	
P15A Ap Applic AAppl	A	14, 14, 14, 14, 14, 14, 14, 14, 14, 14,	14. 44. 44. 44 Ap App A Ap	plicat n ⁷ Appli	Applic:Application	Ap <mark>r</mark>	ApplicApplication	Ap	ApplicAApplication	Applice	Application	AA AA AA
P16A, Ap Applica AAppl	W	14. 44. 44. 44. 44. 44. 44. 44. 44.	44, 44, 44, 44 Ap App A Ap	plicat n ⁷ Appli	Applica Application	Ap	Applic:Application	App ,	ApplicAApplication	MApplic Application	Application	AA AA AA
P17A Ap Applica AAppl	ida - FE,	14, 14, 14, 14, 14, 14, 14, 14, 14, 14,	. #4 #4 #4: #ApiApp A Api	plicat n ⁷ Appli	- Applicat Application	n A <mark>.</mark>	ApplicatApplication	ά n	Applic:AIApplication	Applic Application	Application	AA AA AA.
P18A. Ap Applicat AAppl	iw , FA, FA, FA, FA, FA, FA, FA, FA, FA, F	14, 14, 14, 14, 14, 14, 14, 14, 14,	44, 44, 44, 44 Api App A Api	plicat n ⁷ Appli	ApplicatiorApplicat	ion	ApplicatiorApplicatio	n i	ApplicaticApplication	Applice: Application	Application	AA AA AA.
P19A <mark></mark> Ap Applical AAppl	iv <mark>a</mark> , f.e., f.e., f.f., f.e.,	44, 44, 44, 44, 44, 44, 44, 44, 44,	44 44 44 44 Ap App A Ap	plicat n Appli	ApplicatiApplicatio	n A <mark>l</mark>	ApplicatApplication		ApplicaAApplication	Applic Application	Application	AA AA AA.
P20A, Ap applica AAppl	A <mark>n 148, 44, 44, 44, 44, 44, 44, 44, 44, 44, </mark>	14, 14, 14, 14, 14, 14, 14, 14, 14, 14,	44 44 44 44 App Al Apl	plicat n 🥻 Appli	ApplicatApplicatior	n Ap <mark>l</mark>	, ApplicaApplication	App	applic:AApplication	Applice Application	Application	AA AA AA.
P21A	iv <mark>:</mark> FF, FF, FF, FF, FF, FF, FF, FF, FF	14, 14, 14, 14, 14, 14, 14, 14, 14, 14,	44 44 44 44 AriApp ArAp	plicat n Appli	ApplicatApplication	n Ap <mark>l</mark>	ApplicatApplication	4 <mark>1</mark> 1 - 1	ApplicaApplication	Applice Application	Application	
P22A	iv,	14, 14, 14, 14, 14, 14, 14, 14, 14, 14,	44 44 44 44 App Alap	plicat n <mark>/</mark> Appli	ApplicatApplication	ר A <mark>ן</mark>	ApplicatApplication	₹₽. I	ApplicaApplication	Applice: Application	Application	AA AA AA.
P23Aa Apapplica AAppl	A <mark>n 144 44,</mark> 44, 44, 44, 44, 44, 44, 44, 44,	14, 14, 14, 14, 14, 14, 14, 14, 14, 14,	44 44 44 44 Ap App A Ap	plicat n /Appli	ApplicalApplication	Ap <mark>r</mark>	Application	At is	Applic:AApplication	Applic Application	Application	, AA, AA, AA,
P24A Ap Applic MArApp	lic <mark>o</mark> <i>FA</i> , FA, FA, FA, FA, FA, FA, FA, FA, FA, FA	14. 14. 14. 14. 14. 14. 14. 14. 14.	44, 44, 44, 44 Ap App A Ap	plicat n ⁷ Appli	ApplicApplication	App	ApplicApplication	Арр	ApplicApplication	Applic 4 Application	Application	AA AA /AA /
P25A Ap Appliv ArApp	lic <u></u> <i>FF. FF. FF. FF. FF. FF. FF. FF. FF. FF</i>	14, 14, 14, 14, 14, 14, 14, 14, 14, 14,		plicat n ⁴ Appli	AppliApplication	App	AppliApplication	Ар	AppliArApplication	MApplic Application	Application	
P26A Ap Applic MArApp	lia. <i>14. 14. 14. 14. 14. 14. 14. 14. 14. 14. </i>	14. 14. 14. 14. 14. 14. 14. 14. 14. 14.	44, 44, 44, 44 Ap Appl A Ap	plicat n ⁴ Appli	AppliApplication	App.	ApplicApplication	Арр	AppliApplication	MApplic Application	Application	
P27A Ap Applic MArApp	lic. <i>1</i> 4. 44, 44, 44, 44, 44, 44, 44, 44, 44, 4	14. 14, 14, 14. 14, 14, 14, 14, 14, 14,	44. 44. 44. 44 Appl A Ap	plicat n ⁴ Appli	ApplicApplication	Apr	ApplicApplication	Ap	ApplicArApplication	Applic Application	Application	
P28A Ap Applic MArApp	lica <mark>,</mark> 66, 66, 66, 66, 66, 66, 66, 66, 66, 6	14, 14, 14, 14, 14, 14, 14, 14, 14, 14,	44, 44, 44, 44Ap Appl A Ap	plicat n ⁷ Appli	ApplicApplication	Apr	ApplicApplication	Арр	ApplicApplication	Applic Application	Application	
P29A Ap Applic MArApp	lio, 64, 64, 66, 66, 64, 64, 64, 64, 64, 64	14, 14, 14, 14, 14, 14, 14, 14, 14, 14,	44, 44, 44, 44 Appl A App	plicat n ⁴ Appli	ApplicApplication	App	ApplicApplication	Ap	ApplicArApplication	Applic Application	Application	
P30A Ap ApplivMArApp	lic <i>i, FE</i> , FE, FF, FF, FF, FF, FE, FE, FE, FF, FF	14. 44, 44, 44, 44, 44, 44, 44, 44, 44,	44, 44, 44, 44 Appl A App	plicat n ⁴ Appli	AppliApplication	App <mark>.</mark>	AppliApplication	Apr	AppliApplication	Applic Application	Application	
P31A Appppli MArApp	lic <i>u, FE</i> , FE, FE, FE, FF, FF, FF, FE, FE, FF, FF	145 445 445 447 447 447 447 445 445	## ## ## #Ap Appl A Ap	plication <mark>/</mark> Appli	AppliApplication	Ap <mark>p</mark> iv	AppliApplication	App	AppliArApplication	Applic Application	Application	
P32A	liu a <i>fa. fa. fa fa,</i> fa, fa, fa, fa, fa, fa, fa, fa,	14, 44, 44, 44, 44, 44, 44, 44, 44, 44,	//	plication <mark>/</mark> Appli	AppliApplication	Apr <mark>MI.</mark>	AppliApplication	Ap: <mark>MFA</mark>	AppliA;Application	PApplic Application	Application	
P33A AprApili/MPArApp	lica <mark>,</mark> 66, 66, 66, 66, 66, 66, 66, 66, 66, 6	145 445 445 445 445 445 445 445 445	##########ApplA <mark>,</mark> Apj	plication Appli	AppliApplication	App <mark>M</mark> P	AppliApplication	April	AppliArApplication	Application	Application	
P34A	liw <mark>,</mark> <i>ff,</i> ff, ff, ff, ff, ff, ff, ff, ff, ff,	145 445 445 445 445 445 445 445 445 44	44, 44, 44, 44Ar Appl A, Ap	Cation Appli	Appli Application	App <mark>M</mark>	AppliApplication	App <mark>M</mark>	AppliApplication	Applic Application	Application	
P35A, AppApplic Contemp	I WAY AAY AAL AAL AAL AAY AAL AAL AAY AAY A	Ahre, Ahren Ahren Ahren Ahren Ahren Ahren Ahren Ahren Ahren	, Harr Harr HA. HAR Applican	plication 🛛 🗛 🗛 🗛 🖌	##ApplicApplication	Арр,М	MApplicApplication	App. <mark></mark>	ApplicApplication	Applic Application	Application	AA"AA"AA"

Communication in RTS for 24 rays

Results

IVERSITY

Speedup of NVIDIA Volta V100 over ->	Singlecore	Full MPI node (32 cores)
RTS	27x	1.8x
TVD	36x	2x
MHD	18x	0.78x
Overall	15x	0.9x

Currently:

- 40% runtime is GPU
- 40% runtime is still CPU
- 20% is data movement
- As we finish porting the rest of the code to GPU, and optimizing the parts discussed today, we expect these results to improve significantly

GPU EDUCATION CENTER

Summary

- Accelerating the most significant routine Radiation Transport, among other routines
 - Exploring how to optimize further
- Profiling and re-profiling a Must
- Using directives enable maintenance of a single source code for both multicore and accelerators
 - Enables *new science*

Rich Loft, Eric Wright, Sunita Chandrasekaran <u>loft@ucar.edu</u>, {efwright, <u>schandra}@udel.edu</u>

