



# AWP-ODC Evolution & Large-scale Earthquake Simulation

































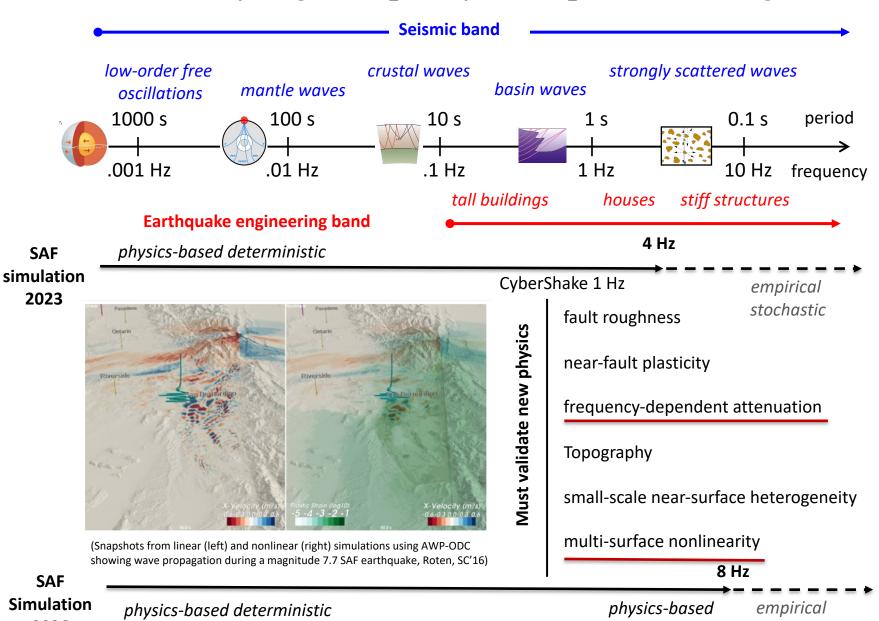








### Why High Frequency Earthquake Modeling



Simulation 2026

stochastic

stochastic



### AWP-ODC

- Started as personal research code (Olsen 1994)
- 3D velocity-stress wave equations

$$\partial_t \mathbf{v} = \frac{1}{\rho} \nabla \cdot \mathbf{\sigma} \qquad \partial_t \mathbf{\sigma} = \lambda (\nabla \cdot \mathbf{v}) \mathbf{I} + \mu (\nabla \mathbf{v} + \nabla \mathbf{v}^{\mathrm{T}})$$

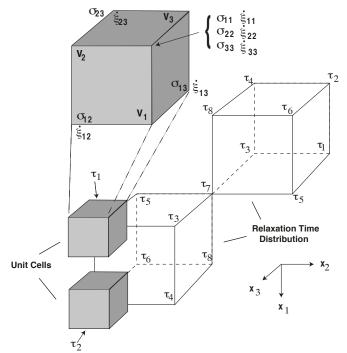
solved by explicit staggered-grid 4th-order FD

**Memory variable formulation of inelastic relaxation** 

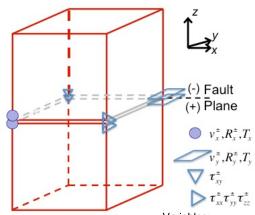
$$\sigma(t) = M_u \left[ \varepsilon(t) - \sum_{i=1}^{N} \varsigma_i(t) \right] \qquad \tau_i \frac{d\varsigma_i(t)}{dt} + \varsigma_i(t) = \lambda_i \frac{\delta M}{M_u} \varepsilon(t)$$
$$Q^{-1}(\omega) \approx \frac{\delta M}{M_u} \sum_{i=1}^{N} \frac{\lambda_i \omega \tau_i}{\omega^2 \tau_i^2 + 1}$$

using coarse-grained representation (Day 1998)

- **Dynamic rupture** by the staggered-grid split-node (SGSN) method (Dalguer and Day 2007)
  - Displacement nodes split at fault surface: explicitly discontinuous displacement & velocity
  - All interactions between sides occur through traction vector at displacement node
- Absorbing boundary conditions by perfectly matched layers (PML) (Marcinkovich and Olsen 2003) and Cerjan et al. (1985)



Inelastic relaxation variables for memoryvariable ODEs in AWP-ODC



Variables:

split-node particle velocities  $V_{i}^{\pm}$ 

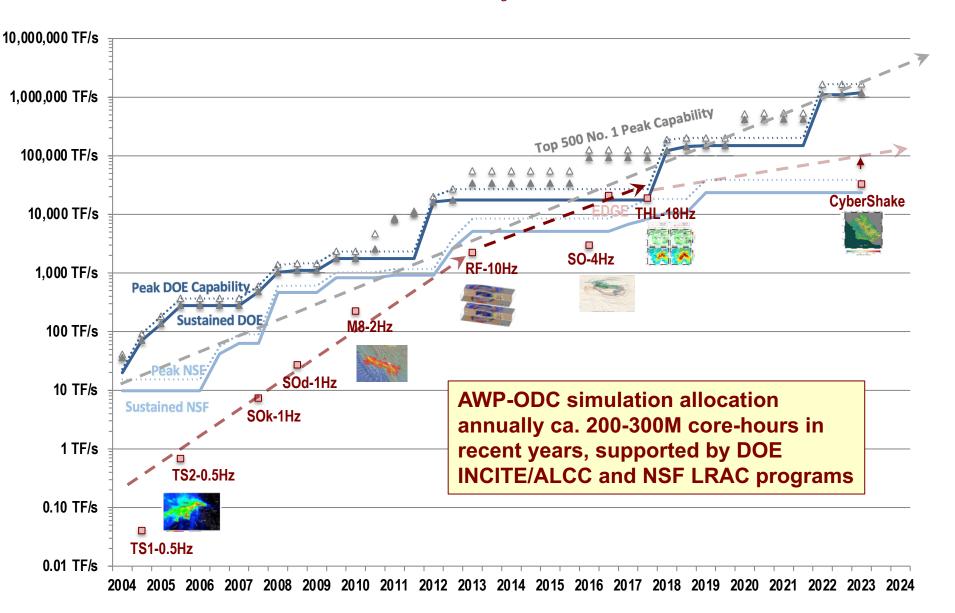
stresses

split-node traction (no iump)

Τį± stress divergence terms



### The Earthquake System Science Challenges at Extreme-Scale Evolution of AWP-ODC



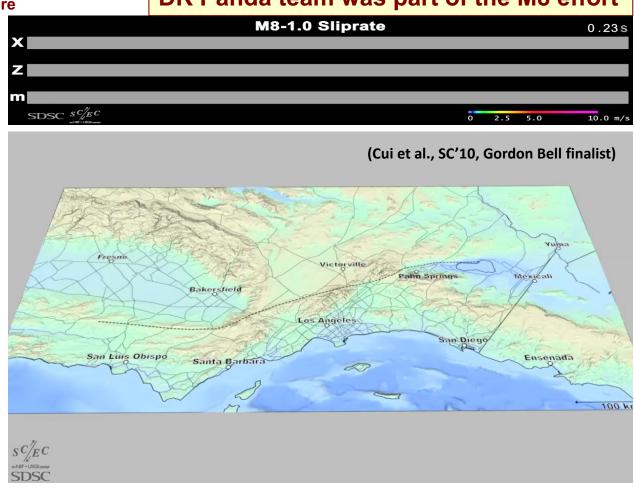


### 0-2 Hz M8 Linear Earthquake Simulation, 2010

- Magnitude 8.0 wall-to-wall scenario, worst-case for southern San Andreas Fault
  - Fault length: 545 km, minimum wavelength: 200 m, NW→SE rupture propagation
- Dynamic rupture simulation performed on Kraken, 7.5 hours using 2160 cores

DK Panda team was part of the M8 effort 881,475 subfaults, 250s of rupture

- Wave propagation simulation performed on Jaguar, 24 hours using 223,074 cores (220 Tflop/s sustained)
  - 436 billion grid points representing SCEC **Community Velocity** Model V4 of dimension 810 x 405 x 85 km (spatial resolution of 40 m)
  - Minimum shear-wave velocity of 400 m/s
  - 368 s of ground motions (160,000 time steps of 0.0023 s) representing seismic frequencies up to 2 Hz

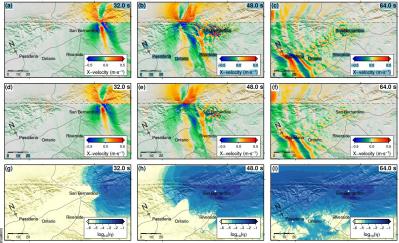


https://www.voutube.com/watch?v=V2Ow0Yuv5co



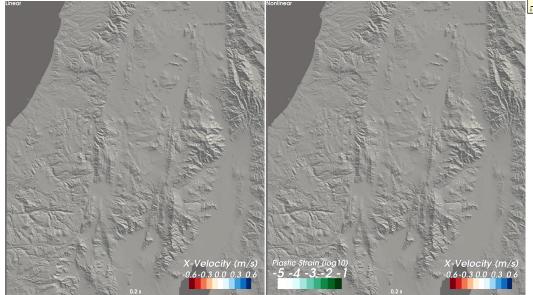
### 0-4 Hz Single-surface J2 Nonlinear ShakeOut Simulation, 2016

- A First 4-Hz nonlinear M7.7 earthquake simulation on the southern San Andreas Fault
- Nonlinear dynamic rupture simulation was conducted using 24,000 CPU-cores on Blue Waters, running 37 hrs
- Nonlinear wave propagation simulation was conducted using 4,200 GPUs on Titan, running 12 hours
- Initially 400% computing time required compared to linear code. With optimized yield factor interpolation, this reduces the computing time from 400% to 165% only



(Roten, et al., SC'16)

- Inside the Whittier Narrows corridor, spectral accelerations at 3 seconds (3s-SAs) are reduced from 1g in the linear case to 0.3-0.6g in the nonlinear case, depending on the choice of reference strain.
- Plastic simulations obtained with a single von Mises yield surface predict 3s-SAs that are higher than those obtained with the multi-surface Iwan model, but lower than the linear values.



ten et al., SC'16) (Roten et al., 2016)

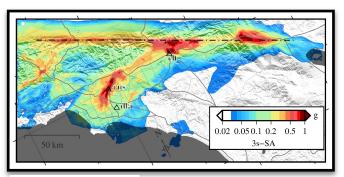


### 0-4 Hz Multi-surface Iwan Nonlinear ShakeOut Simulation, 2023

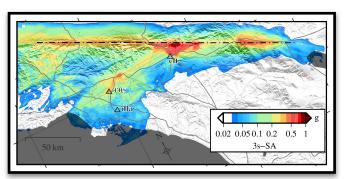
- A multi-surface Iwan type plasticity model in AWP-CPU, verified against the established codes for 1D and 2D SH-wave benchmarks, has been applied to predict the impact of realistic soil nonlinearity on long-period surface waves during large earthquakes on the southern San Andres fault
- While ShakeOut simulations with a single yield surface reduces long period ground motion amplitudes by ~25% inside a wave guide in greater LA, Iwan nonlinearity further reduces the values by a factor of two
- Computational requirements with Iwan model is 20-30x more expensive, and memory use 5-13x more compared to linear solution
- Run 22.5 hrs using 7,680 TACC Frontera nodes



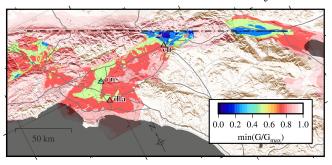
#### Linear



lwan (Darendeli)



Max. shear modulus reduction at the surface



(Roten et al., BSSA, 2023, accepted)

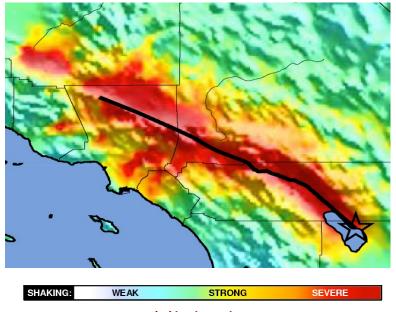


### The ShakeOut Scenario

### M7.8 Earthquake on Southern San Andreas Fault

#### **Scenario Results**

- M7.8 mainshock
  - Broadband ground motion simulation (0-10 Hz)
- Large aftershocks
   M7.2, M7.0, M6.0, M5.7...
- 10,000-100,000 landslides
- 1,600 fire ignitions
- \$213 billion in direct economic losses
  - 300,000 buildings significantly damaged
  - Widespread infrastructure damage
  - 270,000 displaced persons
  - 50,000 injuries
- 1,800 deaths
- Long recovery time



shaking intensity

### Great Southern California ShakeOut

November 13, 2008

### Waveguide amplification in LA Basin

- Caused by string of contiguous sedimentary basins (Olsen et al, 2006, 2009)
- ShakeOut scenario predict strong long-period ground motions in Los Angeles region
- Hazard to pre-Northridge high-rise buildings
- All these approaches assume a linear stress-strain relationship in the fault damage zone and shallow sediments
- Simulations with DP-plasticity predict 30-70% lower ground motions than linear solutions (Roten et al., 2014, 2017)

#### **Exercise Results**

- Largest emergency response exercise in US history, 45M people worldwide participating in 2022
  - Golden Guardian exercise
  - Public events involving multimillion registered participants
- Demonstrated that existing disaster plans are inadequate for an event of this scale
  - Motivated reformulation of system preparedness and emergency response
  - Scientific basis for the LA Seismic Safety Task Force report,
     ilience by Design

Eth New York Eines

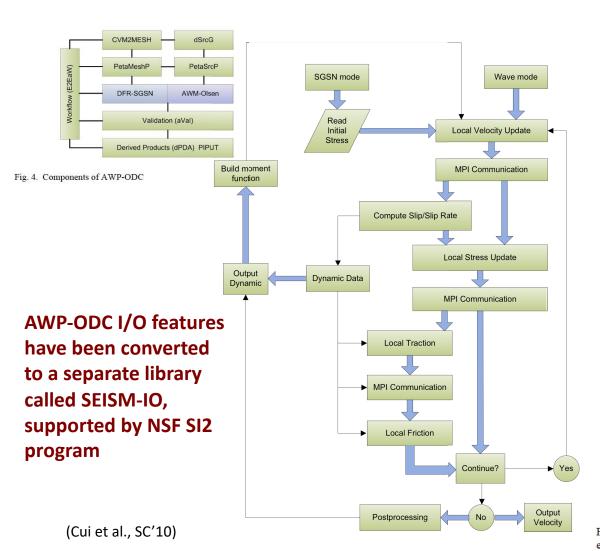
A Seismic Change in Predicting How
by Thomas Fuler

LOS ANGELES — In their quest to make tail buildings safe during earth
are for decades relied on calculations that represent the tremors and
sufficantly underestimate the severity of shaking that buildings of see

Likely to underso durine earthomakes



# Porting to Various CPU Architectures - 2010



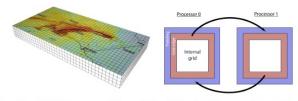


Fig. 5: (left) Decomposition of the M8 simulation region with 810 km long, 405 km wide and 85 km deep; (right) communication between neighboring subgrids

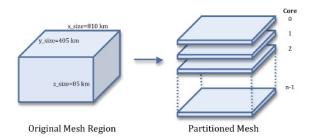


Fig. 7. The 3-D mesh region is partitioned into slices along the z-axis. Each slice is assigned to a core in the MPI job, and each core queries the underlying CVM for the points in its slice only.

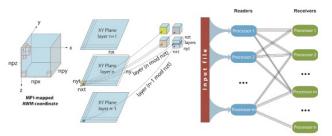


Fig. 9. (left) Cubes and (center) planes for contiguous burst reading and efficient data distributing; (right) high performance I/O with data redistribution



# Porting to GPUs – 2012

- Two-layer 3D domain decomposition on CPU-GPU based heterogeneous supercomputers
  - first step X&Y decomposition for CPUs
  - second step Y&Z decomposition for GPU SMs

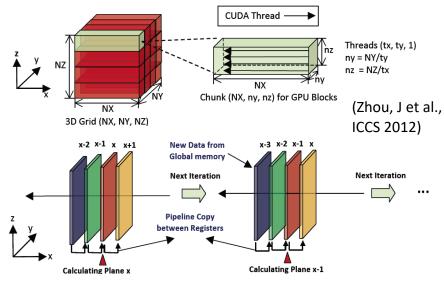
b) Node with multiple a) 3D Domain Decomposition **CPU** cores CPU CPU NY Core Core CPU CPU Core Core CPU Core Core NZ CPU CPU Core Core CPU Memory NX **PCIe** Bus SM SM SM SM Sub-Domain Streaming blocks Multiprocessors mapped to (SM) → CUDA Threads: They compute SM SM NZ SM SM SM decomposed 3D Sub-SM SM Domain blocks c) GPU (in a Node) with multiple SM d) 3D Sub-Domain Decomposition

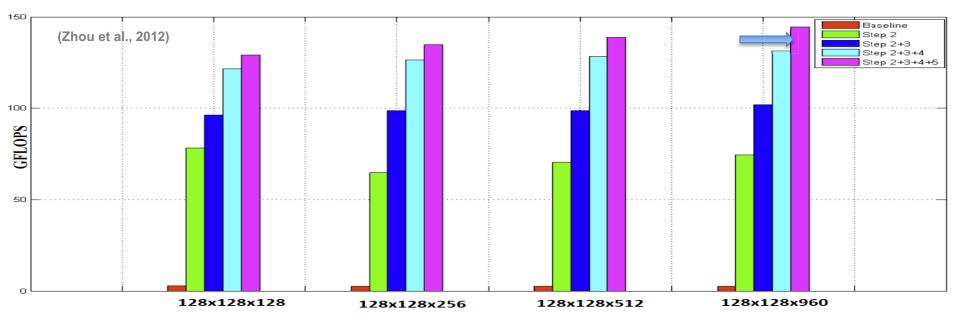
(Zhou et al., ICCS 2012, Cui et al., SC'13)



# Single-GPU Optimizations - 2012

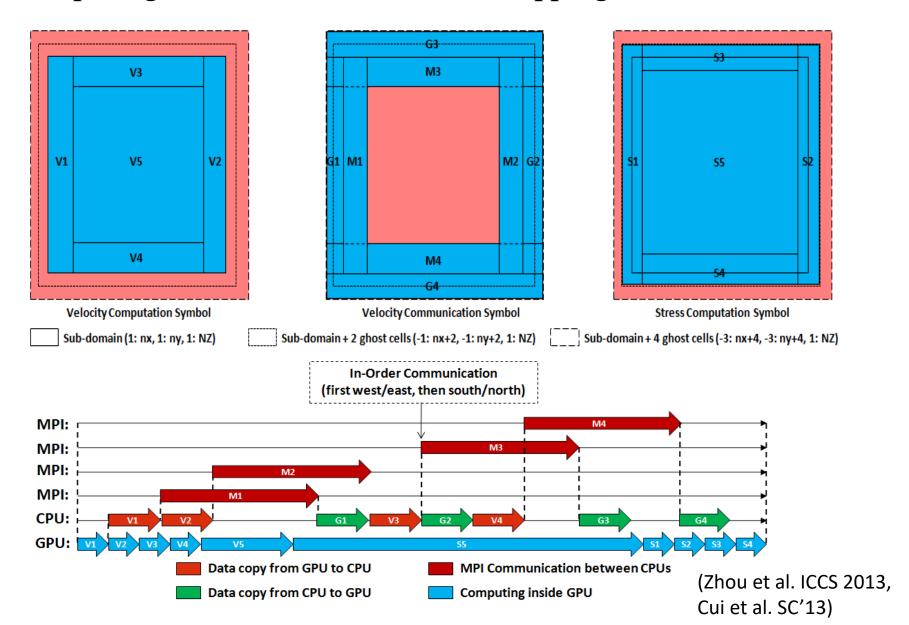
- ✓ Step 2: GPU 2D Decomposition in y/z vs x/y
- ✓ Step-3: Global memory Optimization
  Global memory coalesced, texture memory for six
  3D constant variables, constant memory for scalar constants
- ✓ Step-4: Register Optimization
  Pipelined register copy to reduce memory access
- ✓ Step-5: L1/L2 cache vs shared memory
  Rely on L1/L2 cache rather on-chip shared
  memory





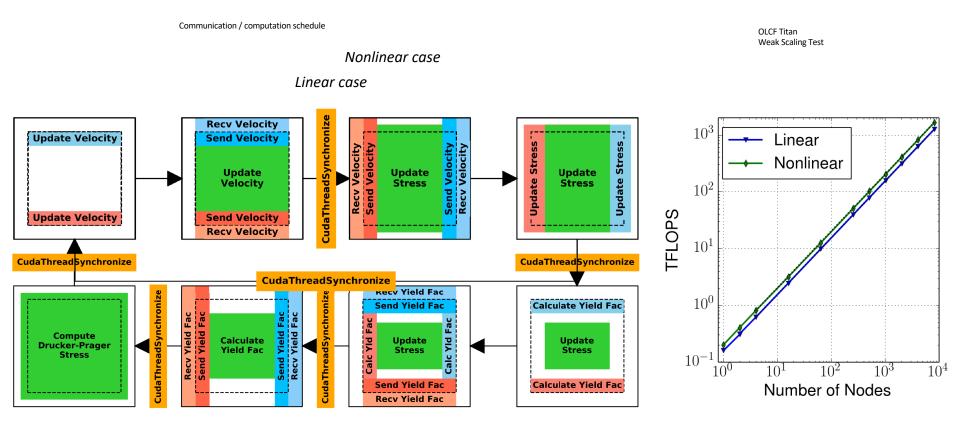


### Computing and Communication Overlapping on GPUs - 2013





# Porting DP-Plasticity on GPUs – 2016





### Porting Iwan Model on CPUs and GPUs – 2021

n = 0

n < nt?

I=0

I < Nspr ?

Yes

No

Yes

Start

l=l+1

Linear stress update

using mu[l], lam[l]

Yield function (von

Mises or DP)

End

### Computational challenges:

 Computationally expensive: separate stress and plasticity update required

for each yield surface

• Memory requirements: each yield surface requires a separate copy of stress tensor  $\tau xx$ ,  $\tau yy$ , ,  $\tau zz$ ,  $\tau xz$ ,  $\tau yz$ ,  $\tau xy$ , Lamé parameters  $\mu$ ,  $\lambda$ , and yield factor r.

MPI communication
 overhead: stress
 tensor and yield
 factor of each yield
 surface needs to be swapped during each time step (reduced scalability)

n=n+1

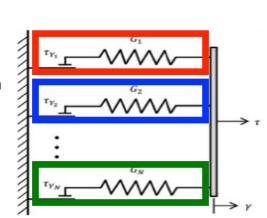
Velocity

update

- Shear modulus reduction reduces max. resolvable frequency
- 10-20x more expensive compared to our 2016 nonlinear simulation which used a simple J2 nonlinear material model, or 20-30x compared to linear solution
- Memory increased by (1 + 0.4\* Nspr) to linear simulation (Nspr = nr of yield surfaces)

#### !wan Concept

 Hysteretic yielding behavior of material represented by a collection of perfectly elasto-plastic spring-slider elements, each element has different constants, shared strain and a fraction of stress, generalized to 3D using a collection of Drucker-Prager yield surfaces



(Roten et al., BSSA 2023, accepted)

No

exceed

tau[l]?

Element

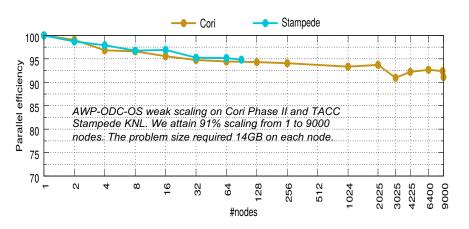
vields

Yes

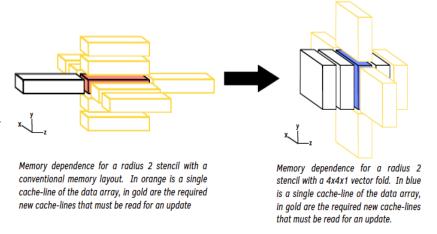


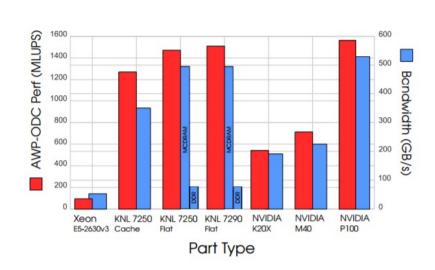
# Porting to Intel Xeon Phi – 2017

- Stencil generation and vector folding through YASK tool: <a href="https://github.com/01org/yask">https://github.com/01org/yask</a>
- Hybrid placement of grids in DDR and MCDRAM
- Normalized cross architecture evaluation in Mega Lattice Updates per Second (MLUPS): Xeon Phi KNL 7290 achieves 2x speedup over NVIDIA K20X, 97% of NVIDIA Tesla P100 performance
- Performance on 9,000 nodes of Cori-II equivalent to performance of over 20,000 K20X GPUs at 100% scaling
- Memory bandwidth accurately predicts performance of architectures (as measured by STREAM and HPCG-SpMv)



(Tobin et al., ISC'17)





Single node performance comparison of AWP-ODC-OS on a variety of architectures. Also displayed is the bandwidth of each architecture, as measured by a STREAM and HPCG-SpMv.



# Porting Discontinuous Mesh on GPUs – 2018

Kim Olsen, SDSU

 $\Rightarrow V_y$ 

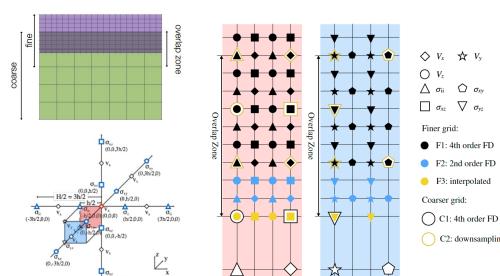
 $\bigcirc$   $\sigma_{\rm rv}$  $\nabla \sigma_{v_7}$ 

F1: 4th order FD

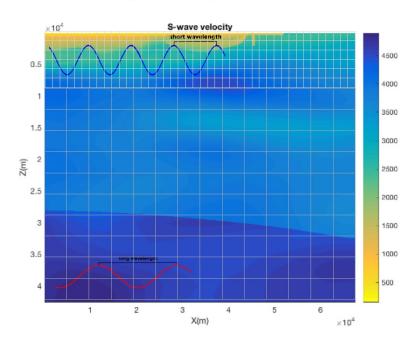
F2: 2nd order FD

C2: downsampling

- Let the interpolation be expressed as: u = W \* U, where U is the field value on the coarse grid, u is the missing point on the fine grid and W is the interpolation operator matrix
- Corresponding downsampling method:  $U' = W^T * u'$ , where u' is the field in the fine grid region, U' is located in the coarse grid, and we set  $W^T$  as downsampling matrix
- Significant performance improvement with respect to a uniform grid solution
  - A factor of 4 achieved for simulating the M9 megathrust earthquake in Cascadia, 650x1000x60 km<sup>3</sup>, 100/300m mesh sizes
  - A factor of 8 achieved for simulating the Mw 5.1 La Habra earthquake up to 4 Hz, using a grid spacing of 20 m in the fine grid and a minimum shear-wave velocity of 500 m/s



wavelength = velocity/frequency



Using a DM with dx<sup>fine</sup> = 100 m in upper 1 km, dx<sup>coarse</sup> = 300 m in bottom 39 km, resulting in 0.28B grids or 72% reduction in grid points

(Nie et al., BSSA 2017, Roten et al., 2018)



# Porting to Topography – 2019

### Christine Goulet, USC

- Topography has been added to AWP-ODC in GPU code, a separate version using curvilinear grids
- Comparable accuracy to the code on a Cartesian grid, with negligible extra memory requirements, longer simulation times due to small timesteps for complex topography
- Perfectly recover a forward simulation using reciprocity a key result needed for CyberShake-related work
- 94% weak scaling efficiency tested up to 1024 GPUs
- ❖ Future plan is to let this curvilinear grid rest on top of layers of Cartesian grids that extend downward with depth

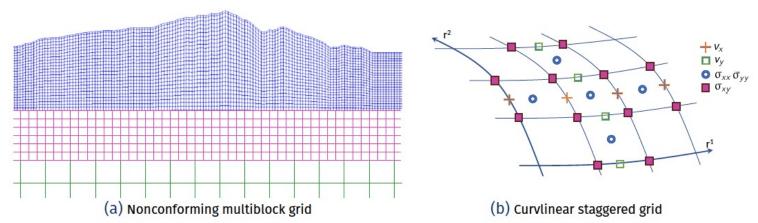


Figure 1: (a) Curvilinear grid, used for discretizing topography, overlaying cartesian grids with decreasing grid resolution with depth. (b) Arrangement of velocity and stresses in a curvilinear staggered grid



# Porting to Microsoft Azure – 2022

### Co-PI: Hari Subramoni

#### **Challenges**

- Digesting the wide breadth of options and configurations
- Higher threshold of initial setup needed
- · Lack of comprehensive forums for debugging errors

- Microsoft Internet2/Azure Accelerator for Research Fall
   2022 program, \$7k credits awarded through Cloudbank
- Future plan is to compare performance with MVAPICH2-AZURE

#### **Benefits**

- Wide flexibility and options of hardware and software allows infrastructure to be tailored to specific workload
- Spin up large VM instances instantly without waiting in a queue/system quotas
- We demonstrated that the AWP performance with a benchmark of ground motion simulation on various GPU based cloud instances, and a comparison of the cloud solution to on-premises bare-metal systems.

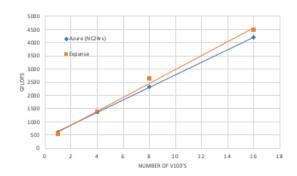


DASHBOARD LEARN ABOUT HELP LOG OUT

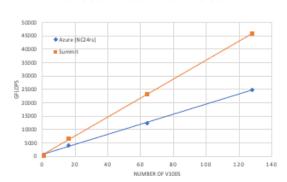
### Accelerating Earthquake Simulation on Microsoft Azure

Specs	Azure (NC24rs)	Expanse	Summit	
GPUs/Node	4 x V100	4 x V100	6 x V100	
СРИ	Xeon E5-2690 v4	Xeon Gold 6248	IBM Power 9	
Memory/Node (GB)	480	384	512	
Compiler:	OpenMPI	OpenMPI	IBM XL Compiler	
File System:	NFS	Lustre	GPFS	
Infinniband (Gbps):	FDR(56)	HDR(200)	EDR(100)	

#### AWP-ODC SCALING ON AZURE VS. EXPANSE



#### AWP-ODC SCALING ON AZURE VS SUMMIT

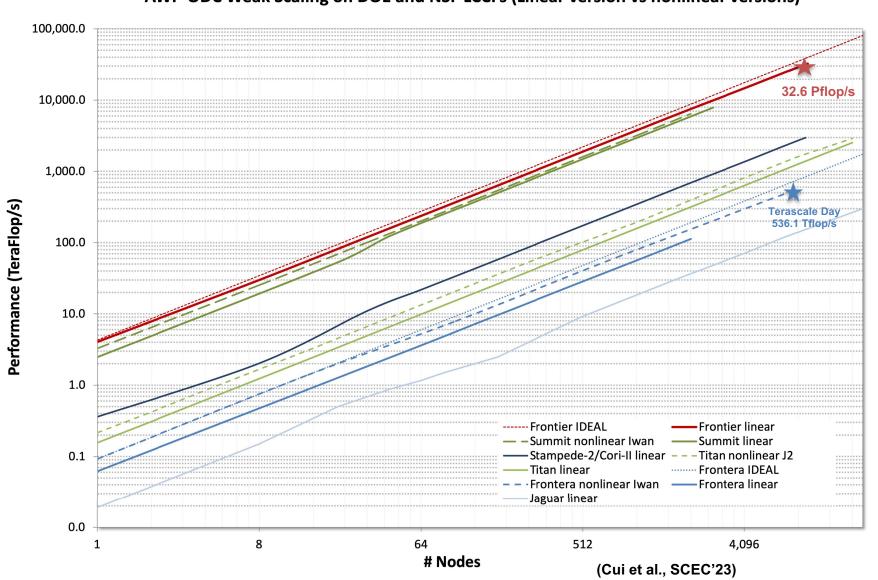


(Palla, SCEC'23)



## Porting CUDA Linear Code to HIP - 2023

**AWP-ODC Weak Scaling on DOE and NSF LCCFs (Linear version vs nonlinear versions)** 



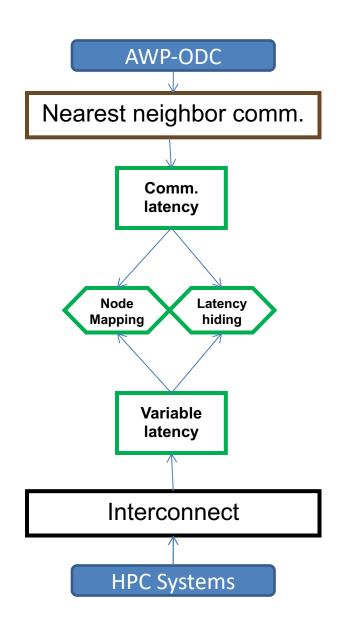


# Accelerate AWP-ODC Performance with MVAPICH









### Performance challenge

- Large variation in communication latencies among neighbors
- System/user memory overhead
- Scalability challenge
  - Increased latency for larger simulation

(Cui et al., SC'10)

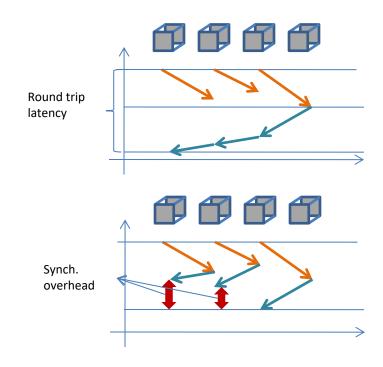
- Asynchronous communication
- Rank re-placement
- Message pre-posting without data reorders
- Computation and communication overlapping, 2-sided and 1-sided





### Asynchronous communication

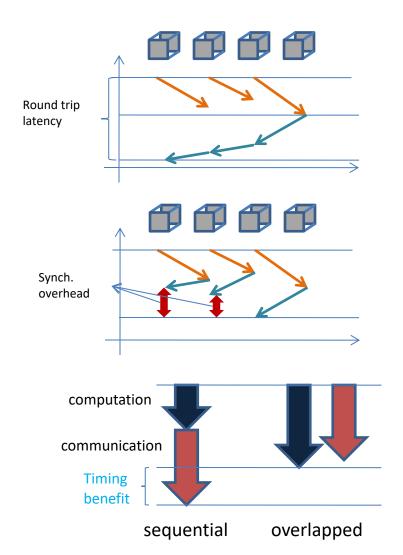
- Significantly reduced latency through local communication
- Reduced system buffer requirement through pre-post receives







- Asynchronous communication
  - Significantly reduced latency through local communication
  - Reduced system buffer requirement through pre-post receives
- Computation/communication overlap
  - Effectively hide computation times
  - Effective when Tcompute\_hide>Tcompute\_overhead







#### **Sreeram Potluri of DK Panda Team**

- Asynchronous communication
  - Significantly reduced latency through local communication
  - Reduced system buffer requirement through pre-post receives
- Computation/communication overlap
  - Effectively hide computation times
  - Effective when Tcompute\_hide>Tcompute\_overhead
  - MPI-1 non-blocking two-sided Communications

#### Velocity Exchange

! recv from south, send to north
n2s(u1, south-mpirank, north-mpirank)
! send to south, recv from north
... repeat for east-west, up-down directions
... repeat for other velocity components v1,w1
wait\_onedirection()
s2nfill(u1, recvbuffer, south-mpirank)

s2nfill(u1, recvbuffer, south-mpirank) n2sfill(u1, recvbuffer, north-mpirank)

s2n(u1,north-mpirank, south-mpirank)

- $\dots\ repeat\ for\ east-west,\ up\text{-}down\ directions$
- ... repeat for other velocity components v1,w1

#### S2N

Copy 2 planes of data from variable to sendbuffer !copy north boundary excluding ghost cells
MPLIsend(sendbuffer, north-mpirank)
MPLIrecv(recvbuffer, south-mpirank)

#### WAIT\_ONEDIRECTION

MPL-Waitall(list of receive requests)

#### S2NFILL

Copy 2 planes of data from recvbuffer to variable ! copy to south ghost cells





#### **Sreeram Potluri of DK Panda Team**

- Asynchronous communication
  - Significantly reduced latency through local communication
  - Reduced system buffer requirement through pre-post receives
- Computation/communication overlap
  - Effectively hide computation times
  - Effective when Tcompute\_hide>Tcompute\_overhead
  - MPI-1 non-blocking two-sided Communications
  - MPI-2 one-sided Communications (on Ranger)

MPI\_Win\_post(group, 0, window) ! pre-posting the window to all neighbors

#### Main Loop in AWM-Olsen

Compute velocity component u
Start exchanging velocity component u
Compute velocity component v
Start exchanging velocity component v
Compute velocity component w
Start exchanging velocity component w
Complete Exchanges of u,v and w
MPI\_Win\_post(group, 0, window)
! For the next iteration

#### Start exchange

MPI\_Win\_start(group, 0, window)
s2n(u1,north-mpirank, south-mpirank)
! recv from south, send to north
n2s(u1, south-mpirank, north-mpirank)
! send to south, recv from north
... repeat for east-west and up-down

#### Complete exchange

MPI\_Win\_wait(window)
MPI\_Win\_complete(window)
s2nfill(u1, window buffer, south-mpirank)
n2sfill(u1, window buffer, north-mpirank)
... repeat for east-west and up-down

#### S2N

Copy 2 planes of data from variable to sendbuffer !copy north boundary excluding ghost cells
MPI\_Put(sendbuffer, north-mpirank)

#### S2NFILL

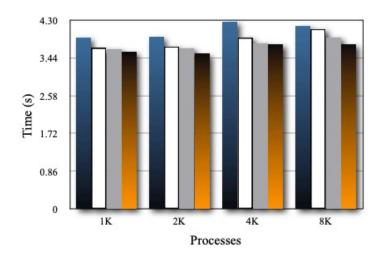
Copy 2 planes of data from window buffer to variable ! copy into south ghost cells

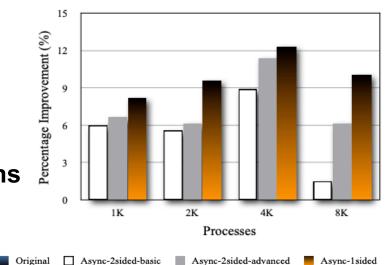




#### **Sreeram Potluri of DK Panda Team**

- Asynchronous communication
  - Significantly reduced latency through local communication
  - Reduced system buffer requirement through pre-post receives
- Computation/communication overlap
  - Effectively hide computation times
  - Effective when Tcompute\_hide>Tcompute\_overhead
  - MPI-1 non-blocking two-sided Communications
  - MPI-2 one-sided Communications (on Ranger)









# Iwan Code Performance on TACC Frontera

module load intel/18.0.5 mvapich2-x/2.3

export MV2\_USE\_MCAST=0

export MV2\_USE\_RDMA\_CM\_MCAST=0

export MV2\_SMP\_EAGERSIZE=28673

export MV2\_SMP\_NUM\_SEND\_BUFFER=8192

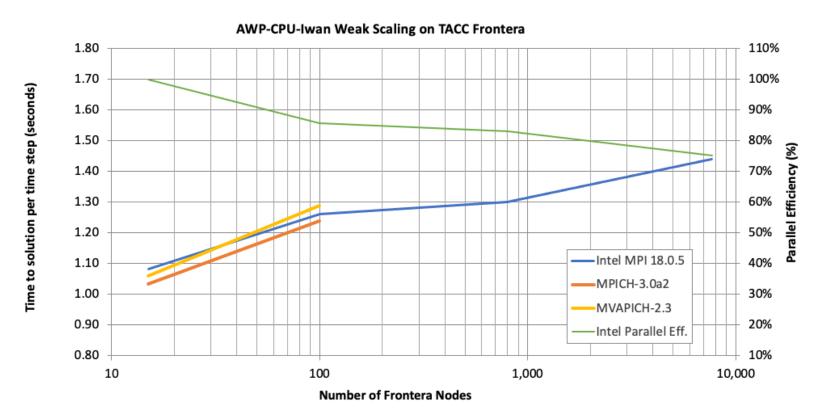
module load intel/18.0.5 mvapich2-x/3.oa2

export MV2\_USE\_MCAST=0

export MV2\_USE\_RDMA\_CM\_MCAST=0

export MV2\_SMP\_EAGERSIZE=28673

export MV2\_SMP\_NUM\_SEND\_BUFFER=8192







# CUDA-aware Support Enhances AWP-ODC Performance

- MVAPICH2 improves performance 20% over OpenMPI on **Expanse**, connected via NVLinks
- MVAPICH2 improves performance 20% over IMPI on Lonestar-6, connected via HDR
- CUDA-aware supported code gains additional 14% in MVAPCICH2/2.37-gdr over MVAPICH-2

Expanse A100s	Teraflop/s	Time (sec/step)		
gcc10.2.0+openmpi4.1.3 (2x2)	2.22	0.0294		
nvhpc21.9 (openmpi4.1.1) (2x2)	2.21	0.0295		
intel19.0.5+mvapich2/2.3.4 (2x2)	2.70	0.0243		
intel19.0.5+mvapich2/2.3.7 (4x2)	3.55	0.0370		
intel19.0.5+mvapich2/2.3.7-gdr (4x2)	4.03	0.0326		

Lonestar 6 A100s	Teraflop/s	Time (sec/step)
gcc11.2.0+impi19.0.9 (2x3)	1.68	0.0585
gcc11.2.0+mvapich2/2.3.7 (2x3)	2.03	0.0488
gcc11.2.0+mvapich2/2.3.7 gdr (2x3)	2.30	0.0399
gcc11.2.0+mvapich2/latest gdr (2x3)	3.15	0.0311





# On-the-fly Compression on GPUs – 2021

### Qinghua Zhou of DK Panda team, IPDPS'21 Best Paper finalist

#### Motivation

- AWP-ODC has significant communication times on large-scale
- Disparity between intra-node and inter-node GPU communication bandwidths that precent efficient scaling

#### **❖** Implementation

- Designed on-the-fly message compression schemes in MVAPICH2-GDR
- Accelerated point-to-point communication performance of transferring large GPU-to-GPU data
- · Compression algorithm for floating-point data, integrated to MVAPICH-GDR
  - MPC: Lossless, high throughput
  - ZFP: lossy, high throughput
- Weak scaling of AWP-ODC on V100 nodes with IB EDR
  - MPC-OPT achieved +18% flops, or -15% runtime
  - ZFP-OPT achieved +35% flops, or -26% runtime

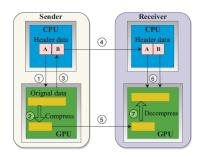


Fig. 4. Data flow of GPU communication with compression. There are seven steps: 1) Launch compression kernel with control parameters 2) Run compression kernel on GPU 3) Returned compressed size 4) Send header data with RTS packet 5) Send compressed GPU data 6) Launch decompression kernel with header data 7) Run decompression kernel to restore the data.

(Q. Zhou et al., IPDPS'21)

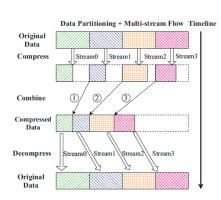
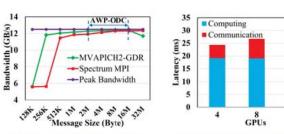


Fig. 7. Data partitioning and multi-stream flow for MPC.

(Q. Zhou et al., IPDPS'21)



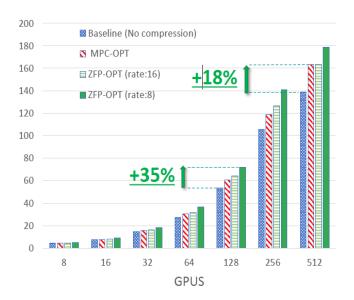
(a) Inter-node D-D Bandwidth

ops (TFLOPS)

(b) AWP-ODC time breakdown

Fig. 2. Motivating Example: production-quality and optimized CUDA-Aware MPI libraries can saturate IB EDR network while the communication time remains a significant bottleneck for HPC applications e.g. AWP-ODC. The message range for AWP-ODC is 2M to 16M as shown in Figure (a).

#### (Q. Zhou et al. IPDPS'21)

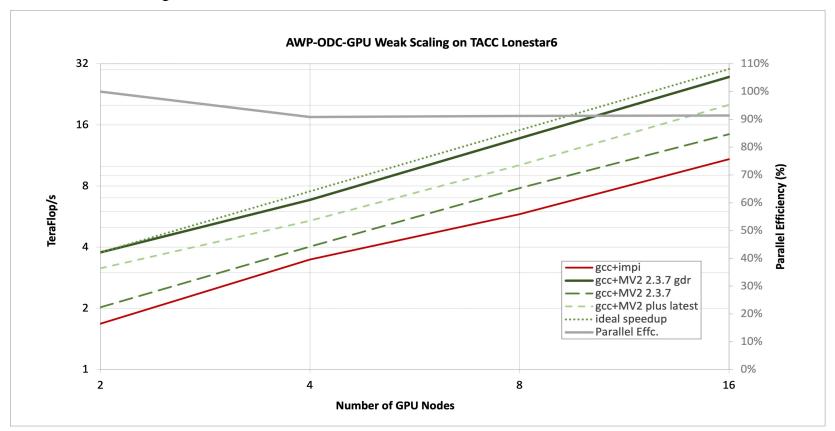


(Q. Zhou et al., IPDPS'21)





# Performance Evaluation on Lonestar-6



- 48%-64% benefits using on-the-fly MPC compression using MPC over GDR
- Combined MVAPICH2-GDR enhancement over IMPI, including both CUDA-aware support and onthe-fly compression, improves application performance by 125%, 97%, 137% and 154% on 2, 4, 8 and 16 nodes, respectively

Lonestar6	mvapich2-2.3.7		mvapich2-2.3.7-gdr		mvapich2-2.3.7-gdr-compresson				
a100	gcc11.2.0			gcc11.2.0		gcc11.2.0			
nodes	Tflop/s	sec/step	parall eff.	Tflop/s	sec/step	parall eff.	Tflop/s	sec/step	parall eff.
2	2.0250	0.0488	100.0%	2.2960	0.0399	100.0%	3.7710	0.0261	100.0%
4	4.0270	0.0494	99.4%	4.5260	0.0436	98.6%	6.8510	0.0288	90.8%
8	7.8250	0.0510	96.6%	9.3250	0.0425	101.5%	13.7560	0.0288	91.2%
16	14.4130	0.1543	89.0%	17.1360	0.0460	93.3%	27.5580	0.0288	91.3%
	impi19.0.9			mvapich2-plus-3.0a2		mvapich2-plus-latest			
li di	gcc11.2.0			gcc11.2.0			gcc11.2.0		
	Tflop/s	sec/step	parall eff.	Tflop/s	sec/step	parall eff.	Tflop/s	sec/step	parall eff.
2	1.6800	0.0585	100.0%	2.391	0.0411	100.0%	3.151	0.0311	100.0%
4	3.4800	0.0572	103.6%	4.579	0.0431	95.8%	5.399	0.0366	85.7%
8	5.8170	0.0686	86.6%	7.796	0.0509	81.5%	10.136	0.0391	80.4%
16	10.8380	0.0737	80.6%	15.214	0.0523	79.5%	20.097	0.0395	79.7%

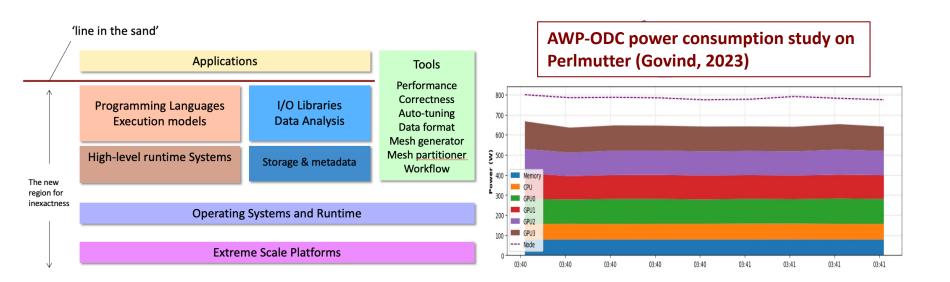


# AWP-ODC software Engineering Challenges and Opportunities



### Challenges for United and Continued Software Development

- Inexact computing is required for reducing energy consumption
- Application level can tolerate a degree accuracy, e.g. discontinuous mesh, error tolerance and precision reduction
- AWP-ODC is highly efficient for regional earthquake simulation and physicsbased seismic hazard analysis





# Summary and Outlook

- AWP-ODC is accelerated with enhanced MVAPICH library on both CPU and GPU architectures
- ❖ We see 154% benefits over IMPI in MVAPICH2-GDR with CUDA-aware support and onthe-fly compression for AWP-ODC on 16 Lonestar6 A100 nodes, future plan is to apply these benefits to Iwan and CyberShake SGT codes
- The Iwan model introduces 20-30x more computation and 5-13x more memory consumption when compared to linear solution, a major challenge for software engineering
- ❖ A joint project with NOWLAB will address load-aware design for MPI asynchronous communication, application-aware neighborhood collective communication, and partitioned point-to-point primitives for efficient communication and cross runtime coordination for MPI+X
- Ongoing NSF CSA project is preparing AWP-ODC for NSF next generation LCCF Horizon to be deployed at TACC – with a hybrid approach using CPUs for dynamic rupture simulation, and GPUs for Iwan-DM wave propagation simulation
- 3D ground motion at 8 Hz or higher is required to realistically capture the full dynamics of a potential Big One on the San Andreas fault



# Acknowledgments

**HPGeoC Team** 

Riverside

Collaborators



**Daniel Roten** 



**Akash Palla** 



**Anish Govind** 



**Philip Maechling** 



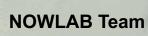
Scott Callaghan



Kim Olsen



Lars Koesterke





**DK Panda** 



Hari Subramoni Sreeram Potluri Qinghua Zhou





**Computing Allocation** 

**OLCF DD, TACC LSCP and CSA, ACCESS Delta, SDSC Expanse, AMD AAC, DOE INCITE & ALCC** 

**Funding** 

NSF LCCF/CSA, NSF CSSI, NSF/USGS SCEC Core, SDSC Core