

Monte Carlo neutron transport in the ECP *Coupled Monte Carlo Neutronics and Fluid Flow Simulation of Small Modular Reactor (ExaSMR) project*



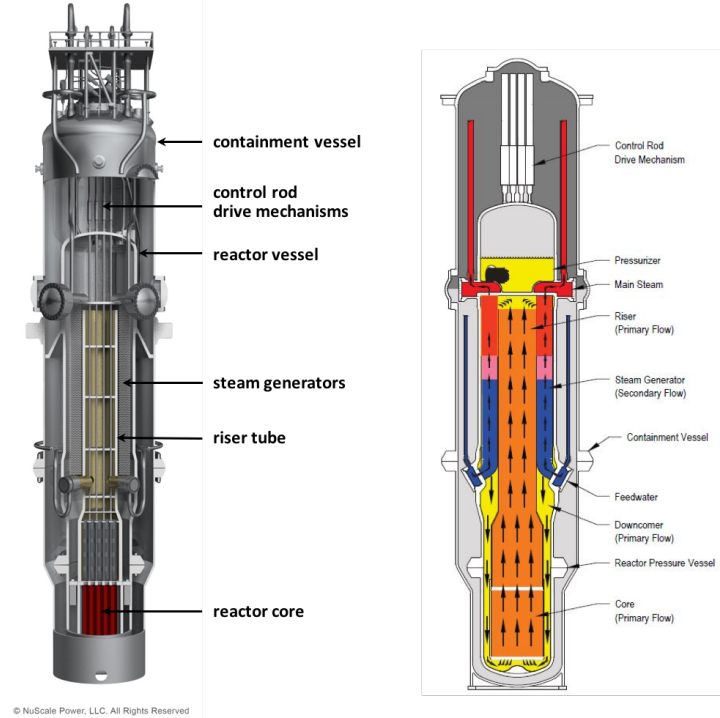
Thomas M. Evans and Steven Hamilton, Oak Ridge National Laboratory

Geant4 R&D Task Force Meeting

April 14, 2020

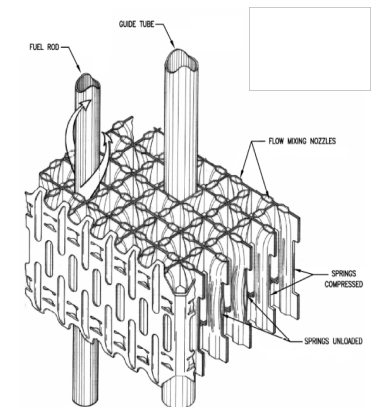
ExaSMR: Modeling and Simulation of Small Modular Reactors

- Small modular nuclear reactors present significant simulation challenges
 - Small size invalidates existing low-order models
 - Natural circulation flow requires high-fidelity fluid flow simulation
- ExaSMR will couple most accurate available methods to perform “virtual experiment” simulations
 - Monte Carlo neutronics
 - Computational Fluid Dynamics (CFD) with turbulence models



© NuScale Power, LLC. All Rights Reserved

Reproduced with permission

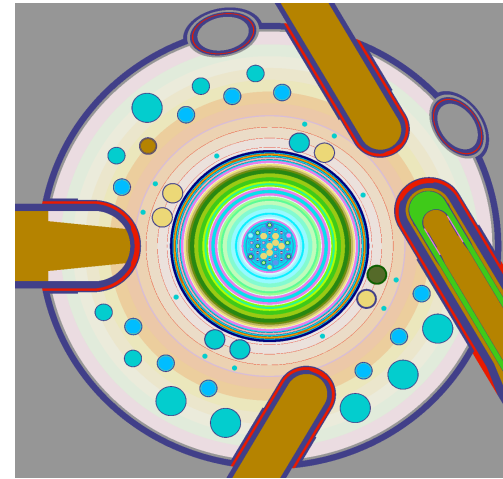


Fuel assembly mixing vane

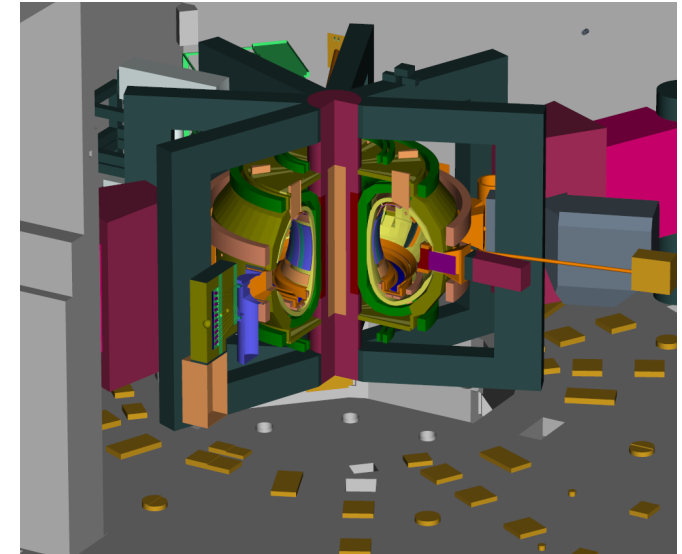
| MC Neutronics | | CFD | |
|--|---|--|---|
| Petascale | Exascale | Petascale | Exascale |
| <ul style="list-style-type: none"> • System-integrated responses • Single physics • Constant temperature • Isotopic depletion on assemblies • Reactor startup | <ul style="list-style-type: none"> • Pin-resolved (and sub-pin) responses • Coupled with T/H • Variable temperatures • Isotopic depletion on full core • Full-cycle modeling | <ul style="list-style-type: none"> • Single fuel assembly • RANS • Within-core flow | <ul style="list-style-type: none"> • Full reactor core • Hybrid LES/RANS • Entire coolant loop |

Transport computational toolkit

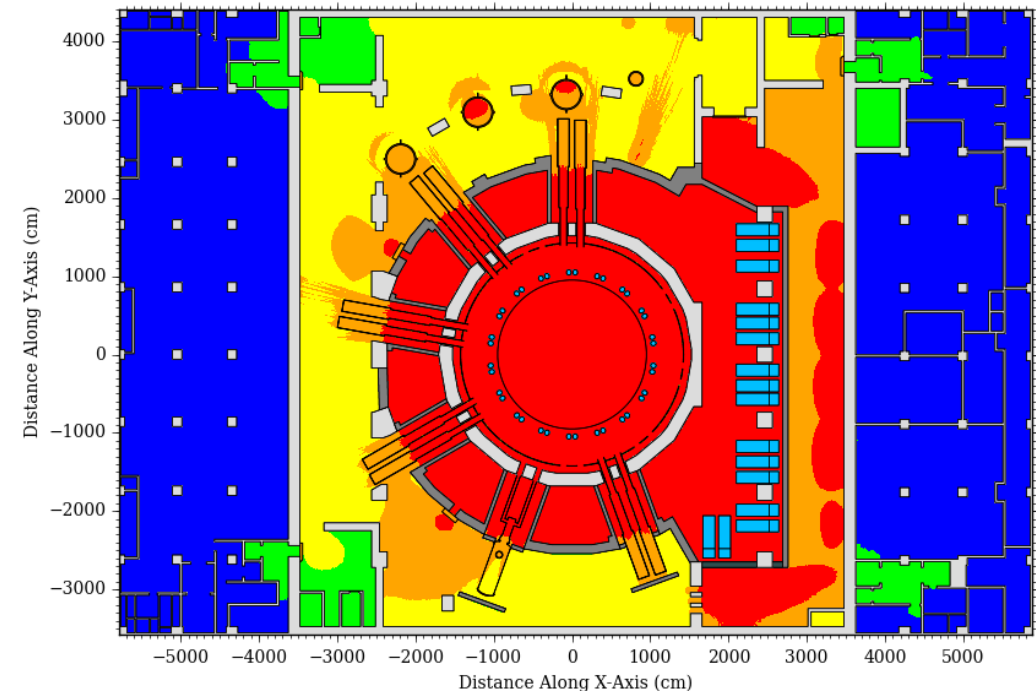
- Exnihilo transport toolkit
 - C++(11/14)/Python/CUDA
 - ~280k lines executable
 - >600k lines including unit-test code
 - 60k lines Python
 - 15k lines in CUDA units (.cu)
- [Shift](#) Monte Carlo neutral particle transport
- [Denovo](#) deterministic solvers
- Hybrid fixed-source (CADIS/FW-CADIS)
- Multiple geometries
 - SCALE solid body general geometry (GG)
 - Reactor ToolKit (RTK)
 - MCNP
 - DagMC-CUBIT CAD
- Physics
 - Multigroup
 - Continuous-energy (CE) tabular
 - Temperature-dependent CE



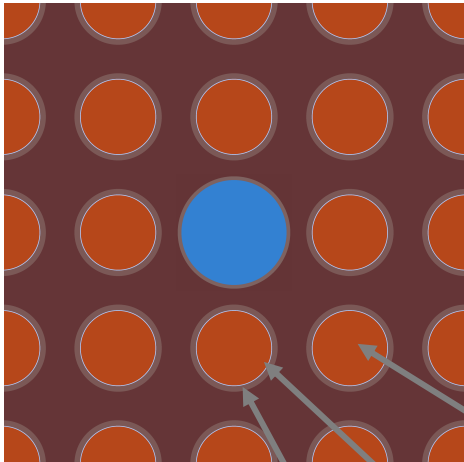
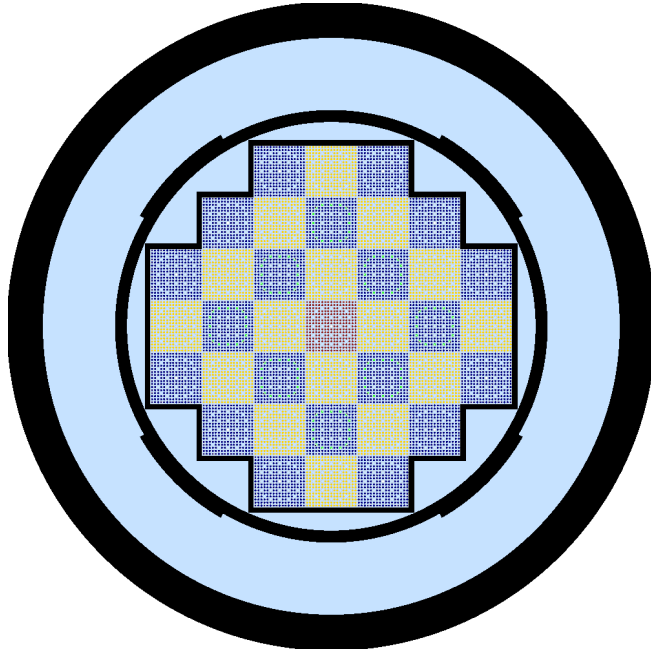
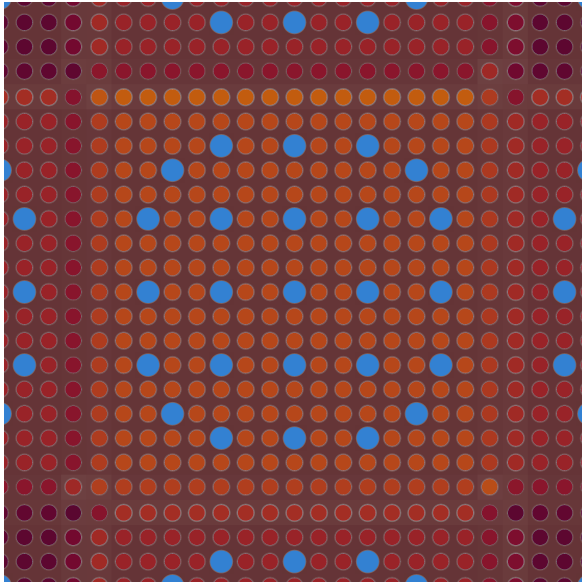
High Flux Isotope Reactor
HFIR



Model of JET facility



Physical Problem Characteristics



$r_f = 0.406$ cm
 $r_g = 0.414$ cm
 $r_c = 0.475$ cm

Pin pitch = 1.26 cm
Assembly pitch = 21.5 cm
Height = 227.56 cm

Fuel (UO_2)

Clad (Zr) Gap (He)

Problem Parameters

• Core Characteristics

- Full core representative SMR model containing 37 assemblies with 17×17 pins per assembly and 264 fuel pins per assembly
- 10^{10} particles per eigenvalue iteration
- Pin-resolved reaction rate with 3 radial tally regions and 50 – 100 axial levels
- O(150) nuclides and O(8) reactions per nuclide in each tally region

• Geometry Size

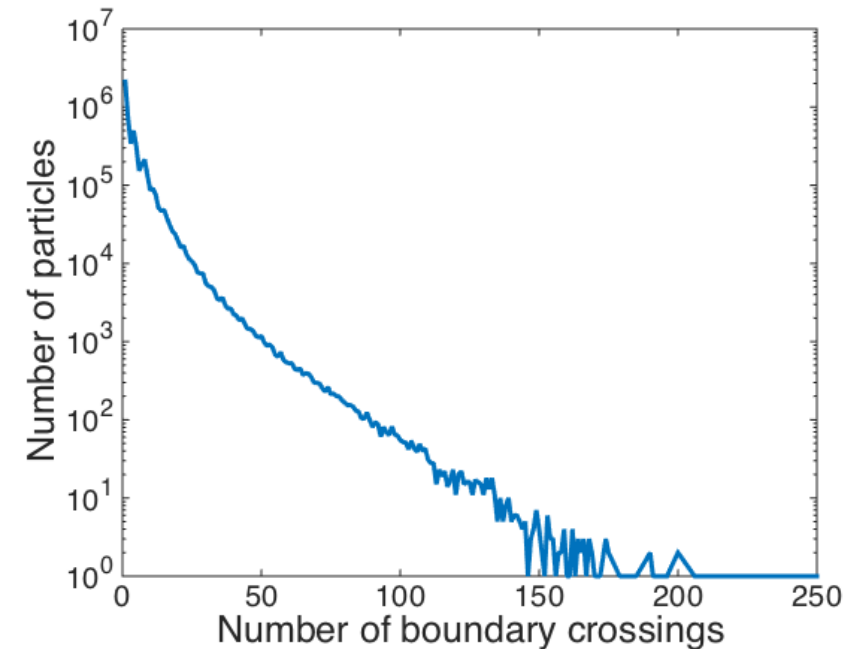
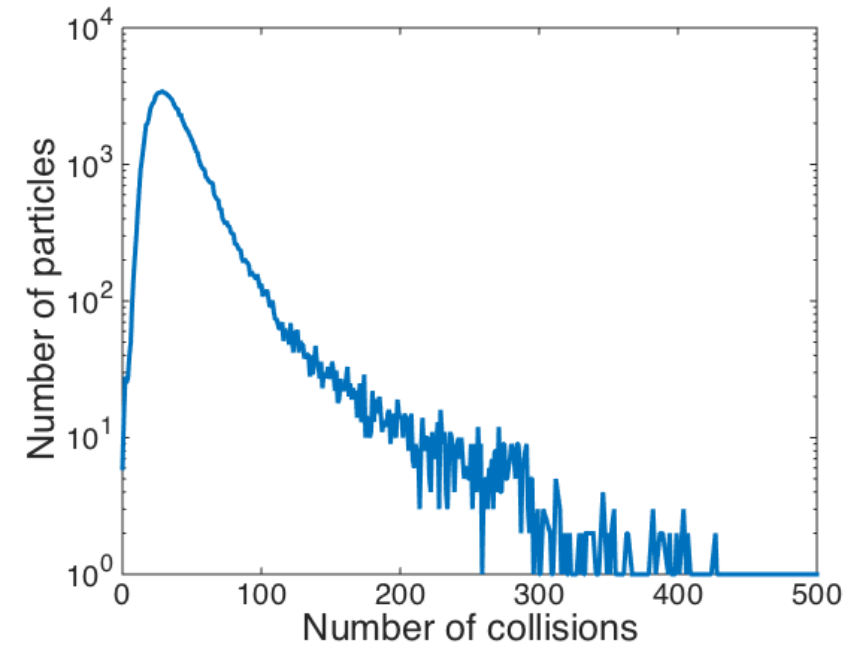
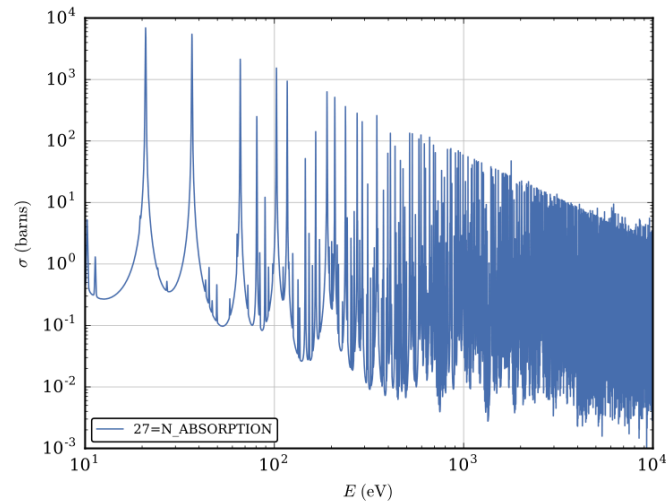
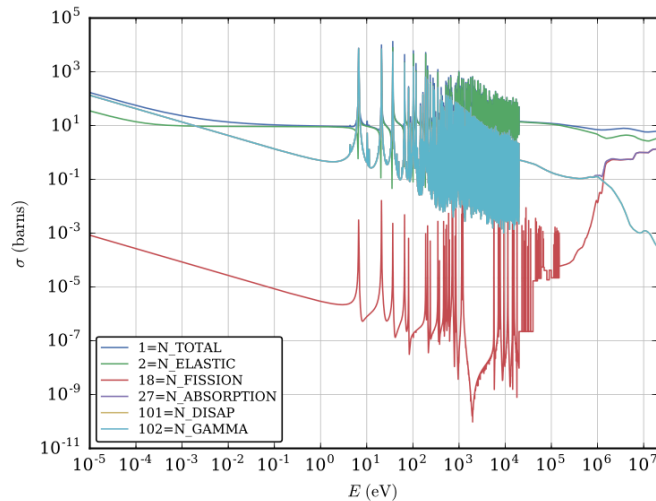
- $N_{cells} = 1.9 \times 10^6 - 8.8 \times 10^6$

• Tally Sizes

- $N_{t,cells} = 4.8 \times 10^5 - 5.9 \times 10^6$
- $N_{t,bins} = 1.5 \times 10^9 - 1.8 \times 10^{10}$

Monte Carlo Neutron Transport Challenges

- MC neutronics is a stochastic method
 - Independent random walks are not readily amenable to SIMT algorithms – on-node concurrency
 - Sampling data is randomly accessed
 - Sampling data is characterized by detailed structure
 - Large variability in transport distributions both within and between particle histories



Developing GPU Continuous Energy Monte Carlo – Intra-Node

- Focus on high-level thread divergence
- Optimize for device **occupancy**
 - Separate geometry and physics kernels to increase occupancy
 - Boundary crossings (geometry)
 - Collision (physics)
- Smaller kernels help address variability in particle transport distributions
- Partition macro cross section calculations between fuel and non-fuel regions – separate kernels for each
- Use of hardware atomics for tallies and direct sort addressing
- Judicious use of *texture* memory
 - `__ldg` on data interpolation bounds

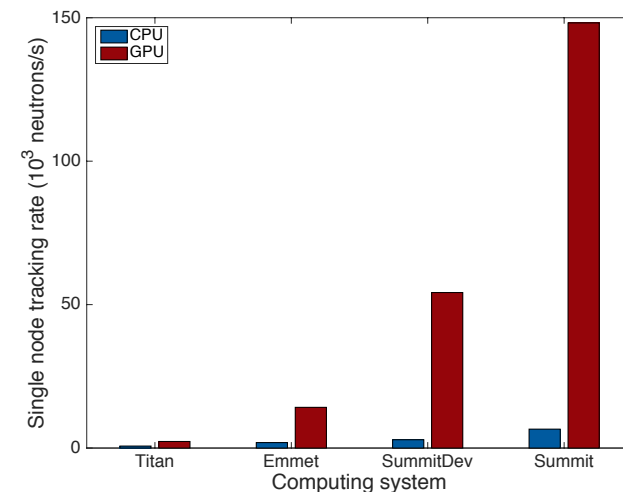
Simple Event-Based Transport Algorithm

```
get vector of source particles
while any particles are alive do
  for each living particle do
    move particle
    dist-to-collision
    dist-to-boundary
    move-to-next
  end for
  for each living particle do
    process particle collision
  end for
  source particles
  sort/consolidate surviving particles
end while
```

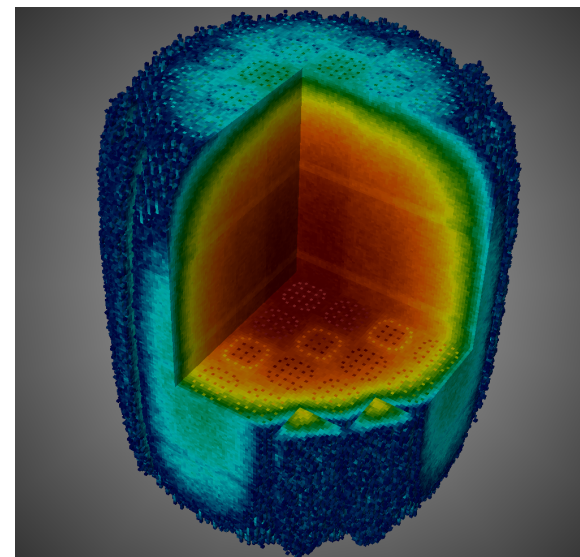
Production continuous-energy Monte Carlo transport solver on GPUs

- Petascale implementation did not use GPU hardware
- Enables three-dimensional, fully-depleted SMR core models simulated using continuous-energy physics and pin-resolved reaction rates with temperature-dependence
- Algorithmic improvements offer 10x speedup relative to initial implementation and nearly 60x per-node speedup over Titan
- Nearly perfect parallel scaling efficiency on ORNL's Summit supercomputer
- GPU algorithm executes more than 20x faster than CPU algorithm on Summit (per full node)
- Paper describes first production MC solver implementation on GPUs

Hamilton, S.P., Evans, T.M., 2019. Continuous-energy Monte Carlo neutron transport on GPUs in the Shift code. *Annals of Nuclear Energy* **128**, 236–247. <https://doi.org/10.1016/j.anucene.2019.01.012>



Increase in particle tracking rate across GPU computing architectures



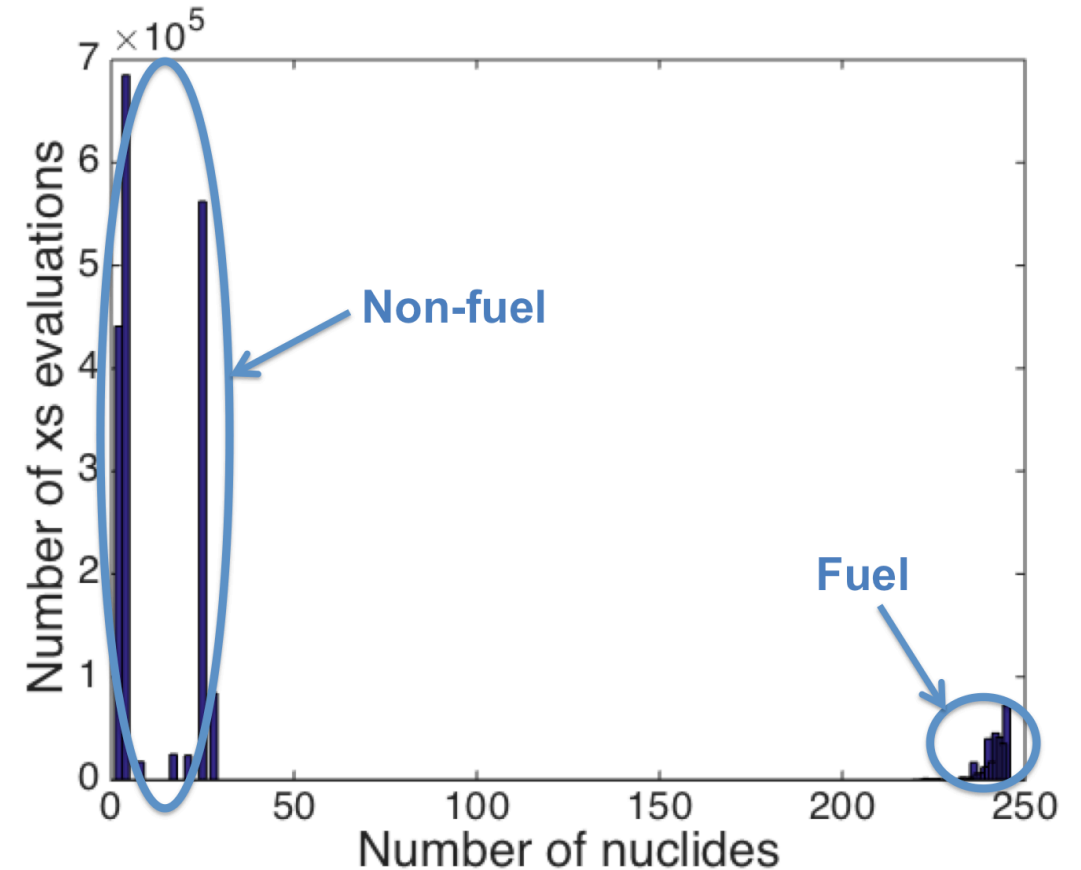
Total reaction rate in SMR core

Cross section calculations

- Computing transport cross sections requires contributions from various constituents

$$\Sigma(E) = \sum_{m=1}^M N_m \sigma_m(E)$$

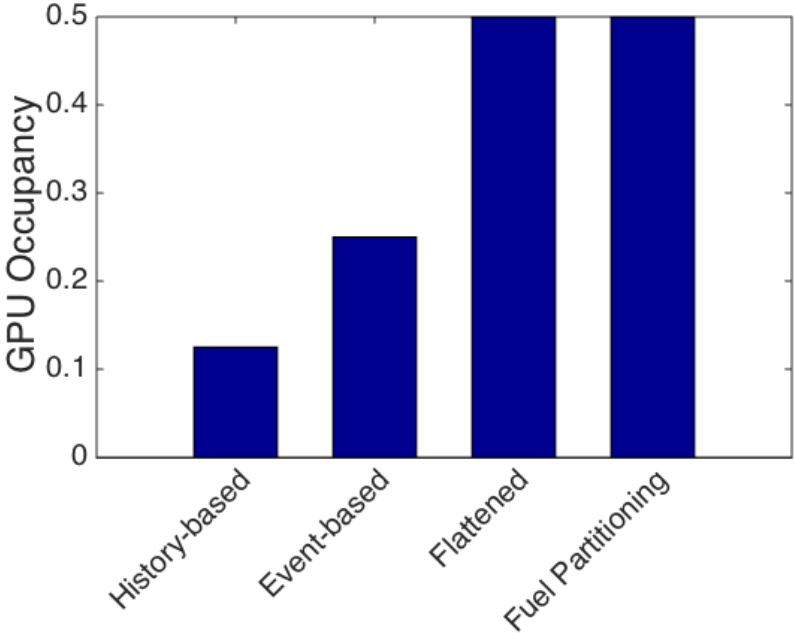
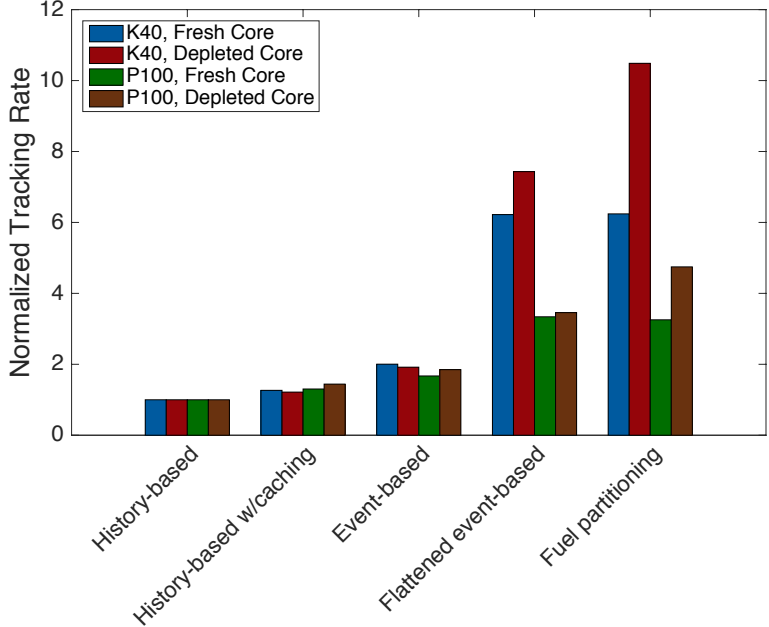
- Fuel compositions contain substantially more nuclides than non-fuel
- Partition mixtures into fuel and non-fuel
 - Evaluate cross sections in separate kernels to reduce divergence



GPU Device Occupancy

- Flattened algorithm allows small, focused kernels
 - Split geometry/physics components to reduce register usage
 - Smaller kernels = higher occupancy

| MC type | Algorithm | Registers | Occupancy |
|-------------------|---------------|-----------|-----------|
| Multigroup | History-based | 85 | 25% |
| | Event-based | 83 | 25% |
| Continuous-Energy | History-based | 168 | 12.5% |
| | Event-based | 62 | 50% |



Effect of varying occupancy

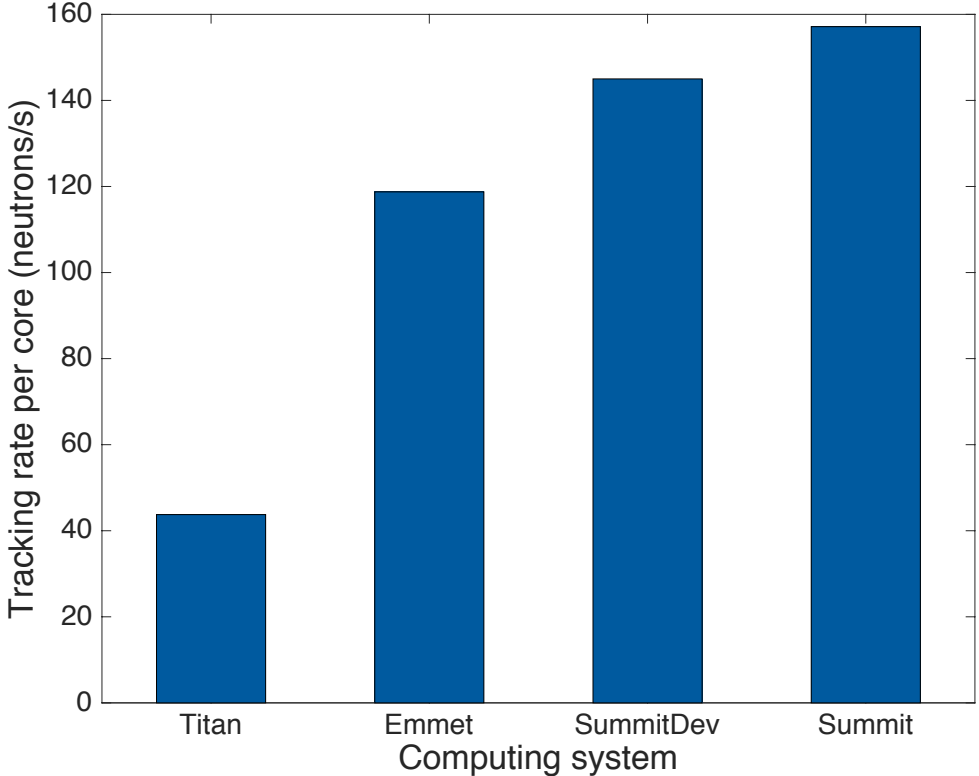
- Artificially limit occupancy by allocating shared memory
 - `kernel<<< grids, blocks, shared_mem>>>(...)`

| Occupancy (%) | Algorithm | | |
|---------------|---------------|-------------|-----------------------|
| | History-based | Event-based | Flattened event-based |
| 12.5 | 3.7 | 3.4 | 8.2 |
| 25.0 | - | 5.8 | 13.3 |
| 37.5 | - | - | 14.5 |
| 50.0 | - | - | 16.9 |
| 62.5* | - | - | 18.0 |

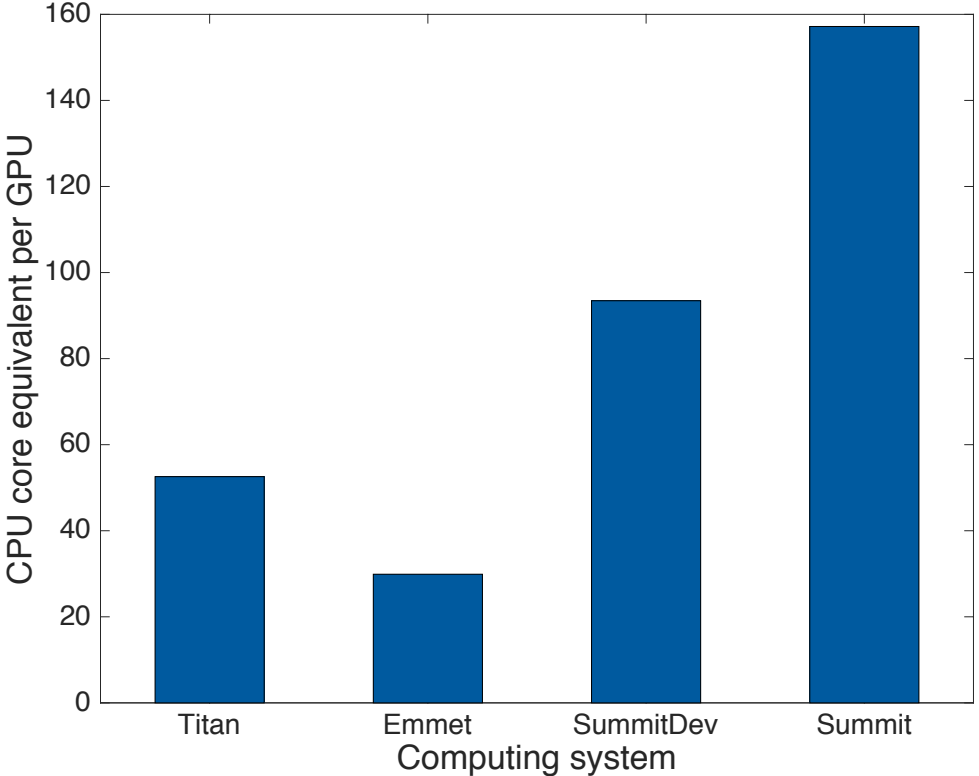
*Only applied to “distance to collision kernel”

CPU v GPU performance

CPU tracking rate per core

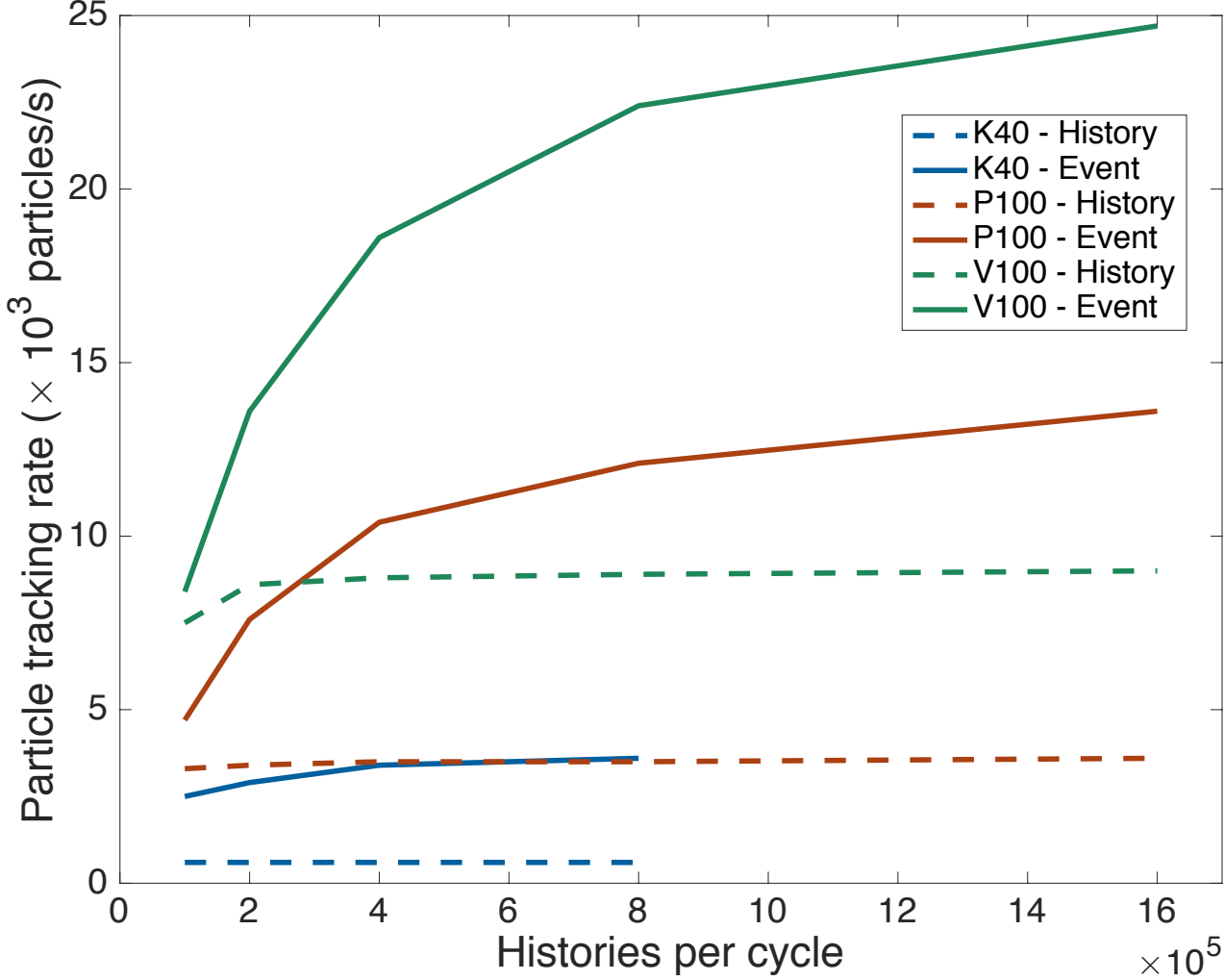


GPU core equivalent



GPU performance increases have outpaced corresponding CPU improvements

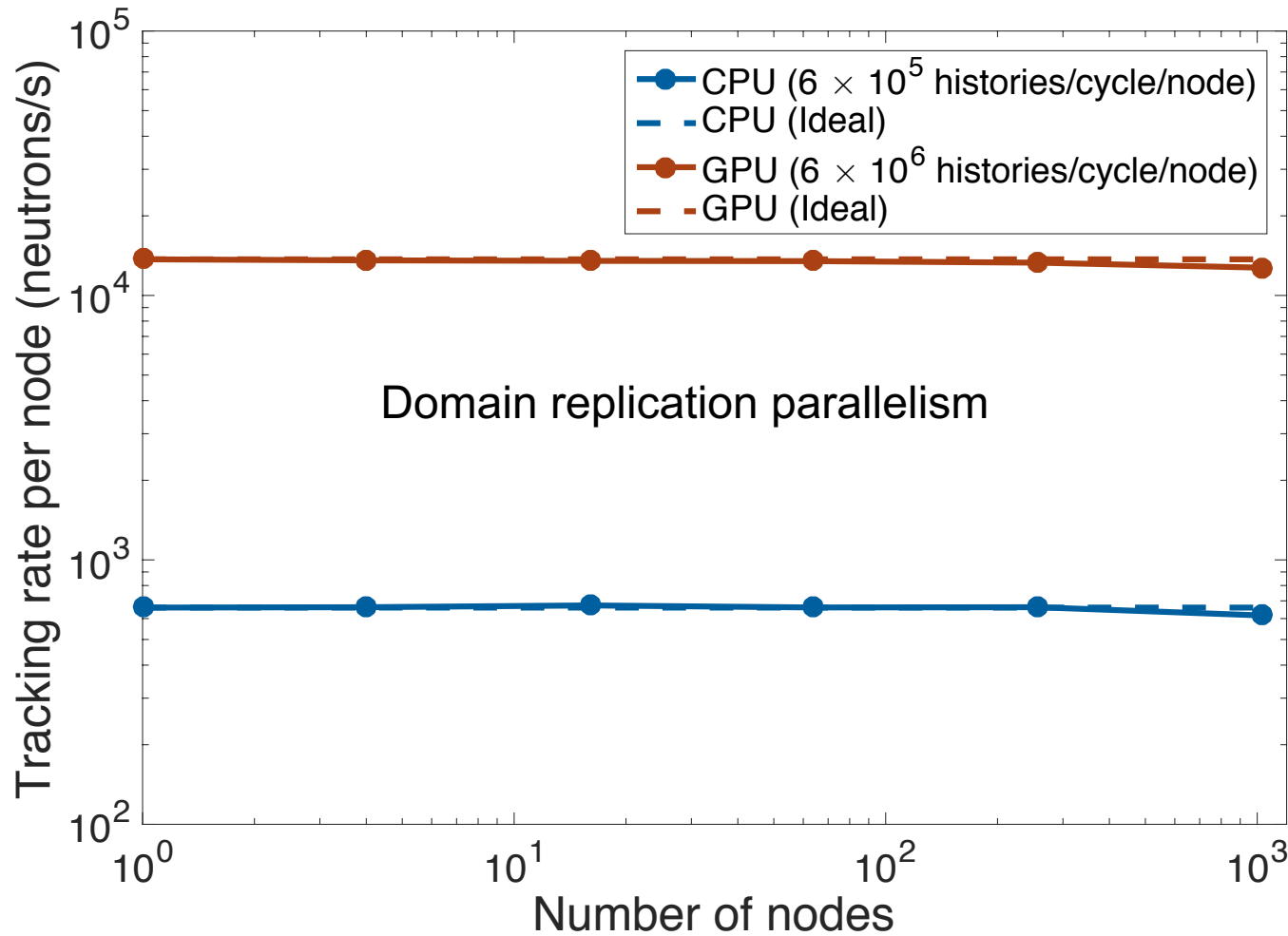
Device saturation



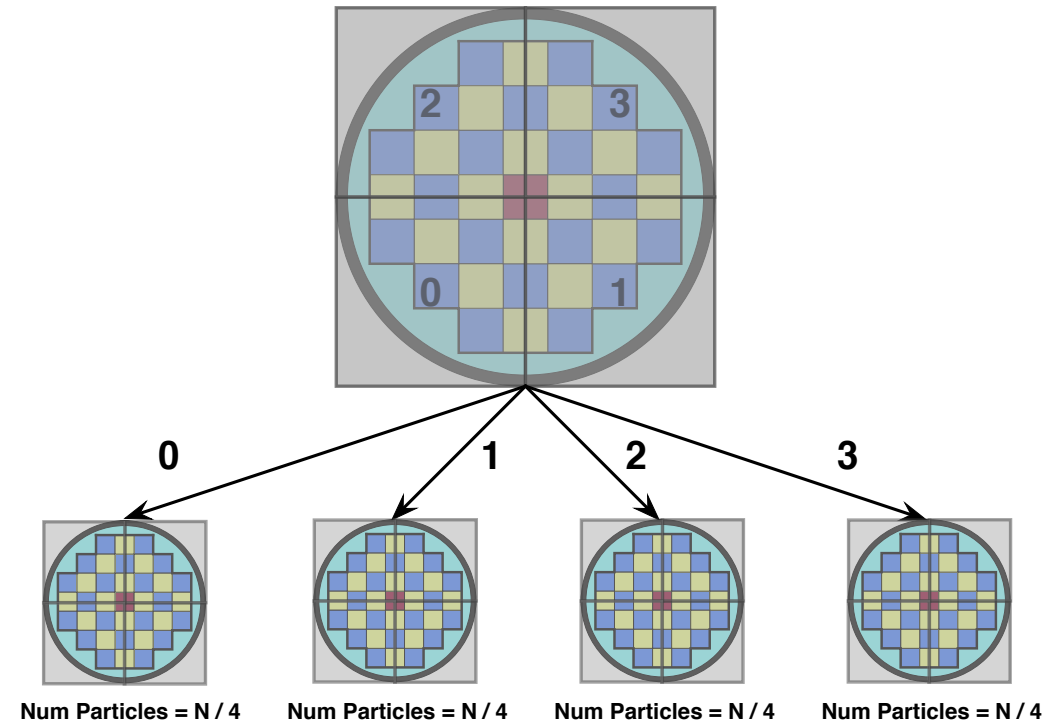
Newest architectures remain unsaturated at 1M particles per GPU

Depleted SMR core

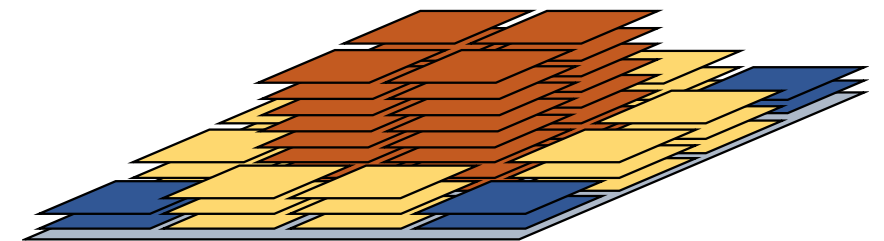
Inter-node Scaling



Weak scaling on Summit – 1 GPU per MPI rank



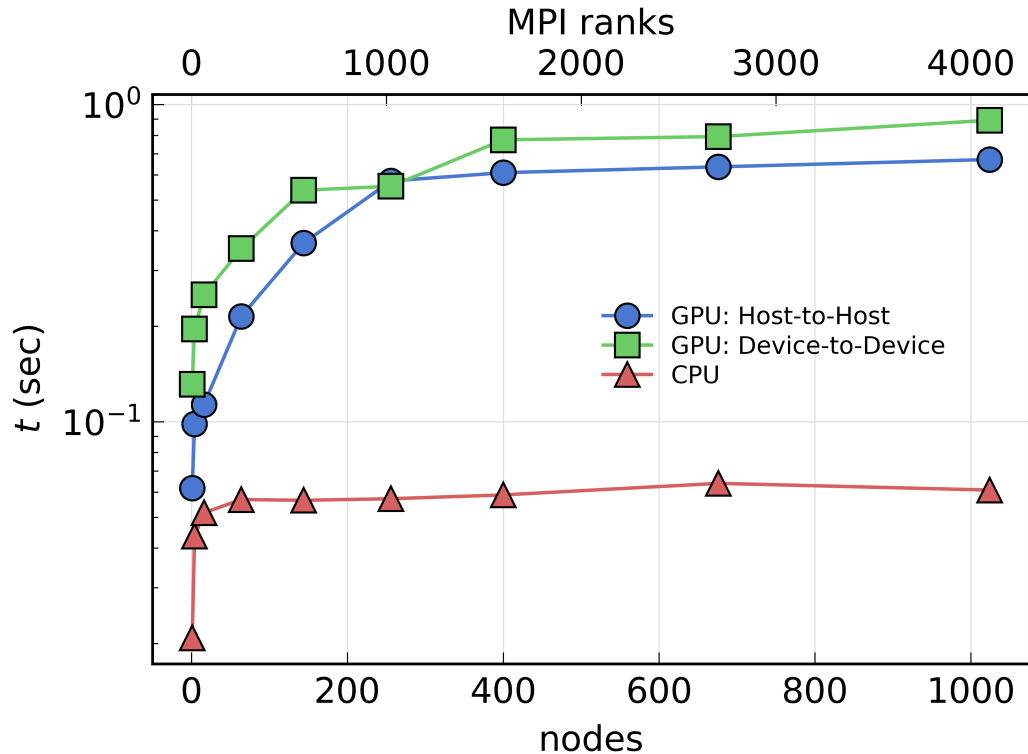
Multi-set domain decomposition topology
(in development – GPU)



Intra-set non-uniform block out to address load balancing

Domain decomposition fits cleanly in event algorithm

2 Resource Sets per Node (2 GPUs per Resource Set)



Investigating MPI-aware CUDA

- Communication device-to-device (bypass NIC)
- Does not currently give same performance as manually moving data
- Next-gen platforms will optimize device-to-device

Event-Based Transport Algorithm

```
while have active particles do
  for iter=1, Ngen do
    if event[SOURCE] then
      source_event(bank)
    if event[MOVE_TO_NEXT] then
      move_event()
    if event[COLLISION] then
      collision_event()
    if event[BUFFER] then
      buffer_event()
      update_inactive()
      process_indices()
    end for
    communicate(bank)
  end while
```

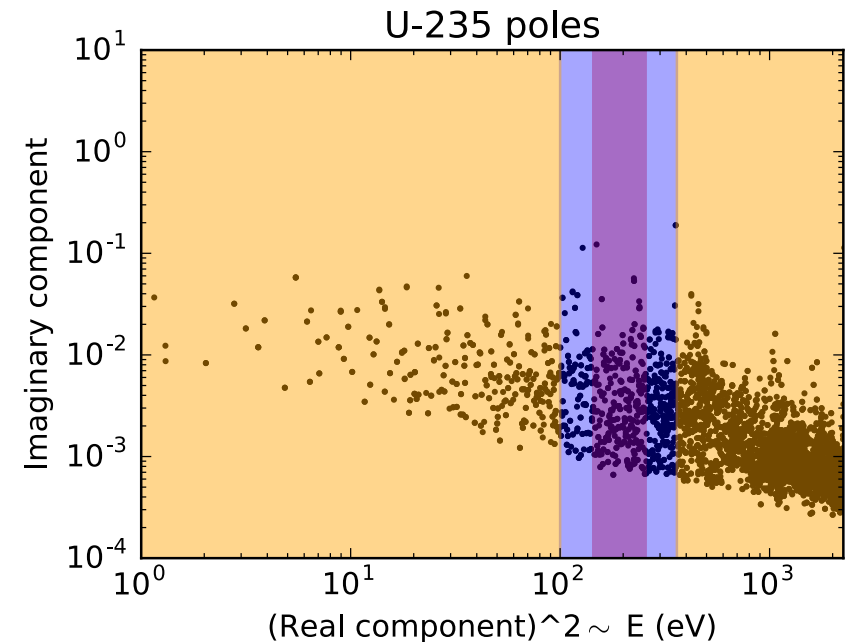
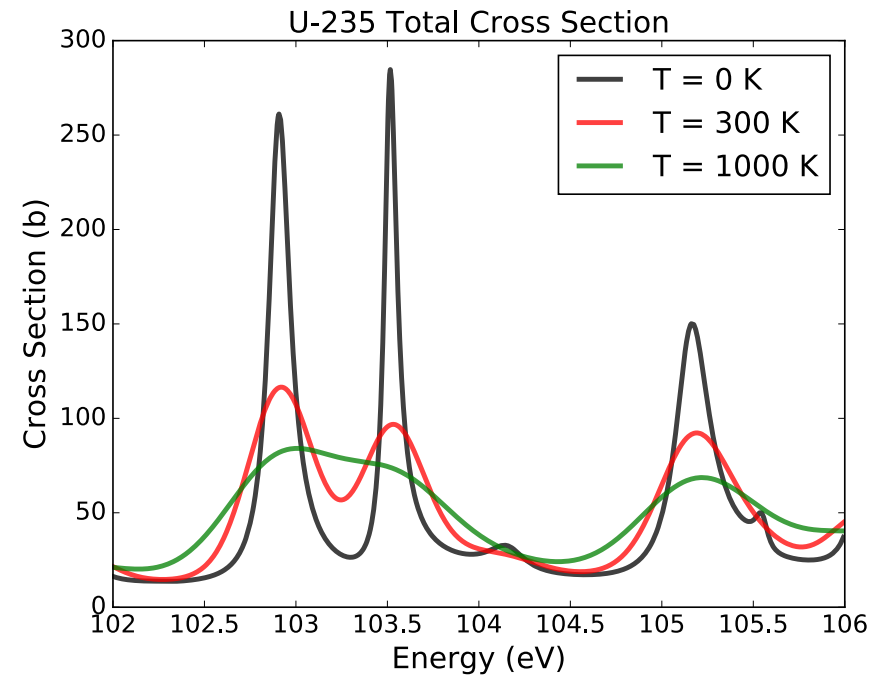
On-the-Fly Doppler Broadening

- Cross section resonances significantly broaden due to thermal motion of nuclei
- The cross section (σ) at any energy (E) and temperature (T) can be expressed as a summation over contributions from poles (p_j) and corresponding residues (r_j):

$$\sigma(E, T) = \frac{1}{E} \sqrt{\frac{A\pi}{k_B T} \sum_j \Re \left[r_j W \left((\sqrt{E} - p_j) \sqrt{A/k_B T} \right) \right]}$$

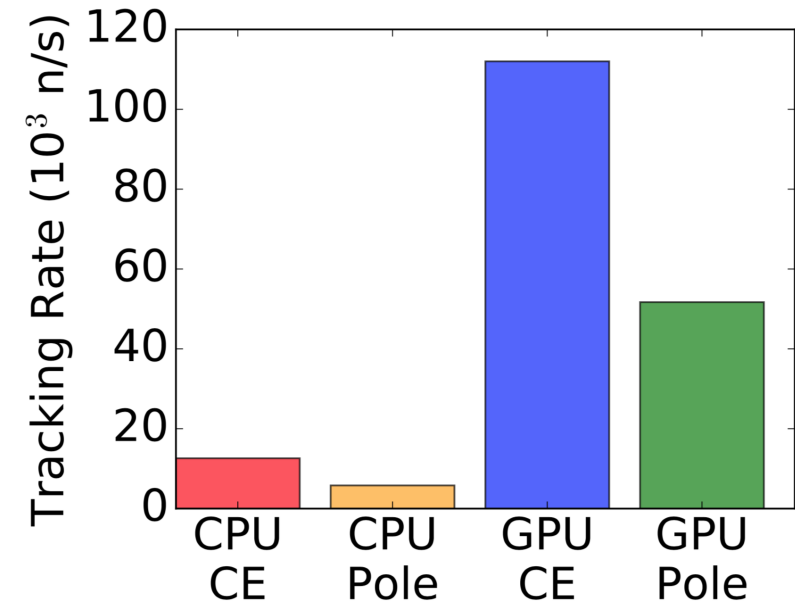
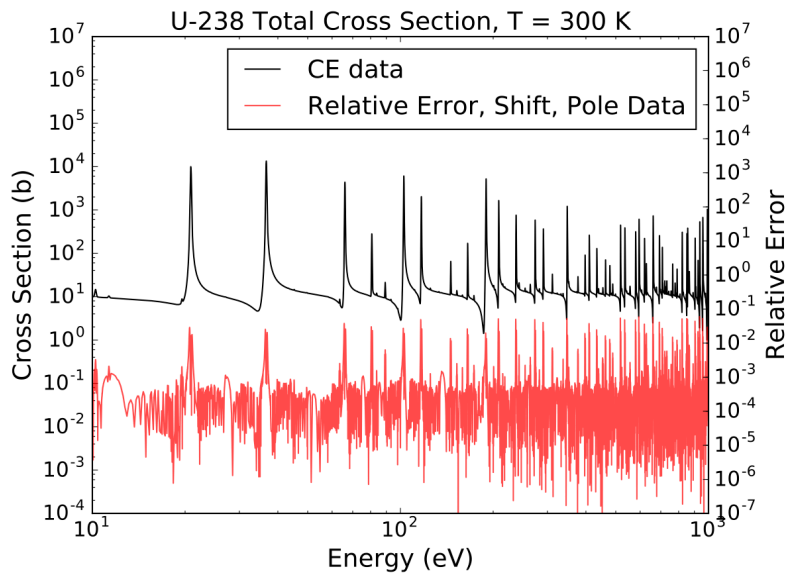
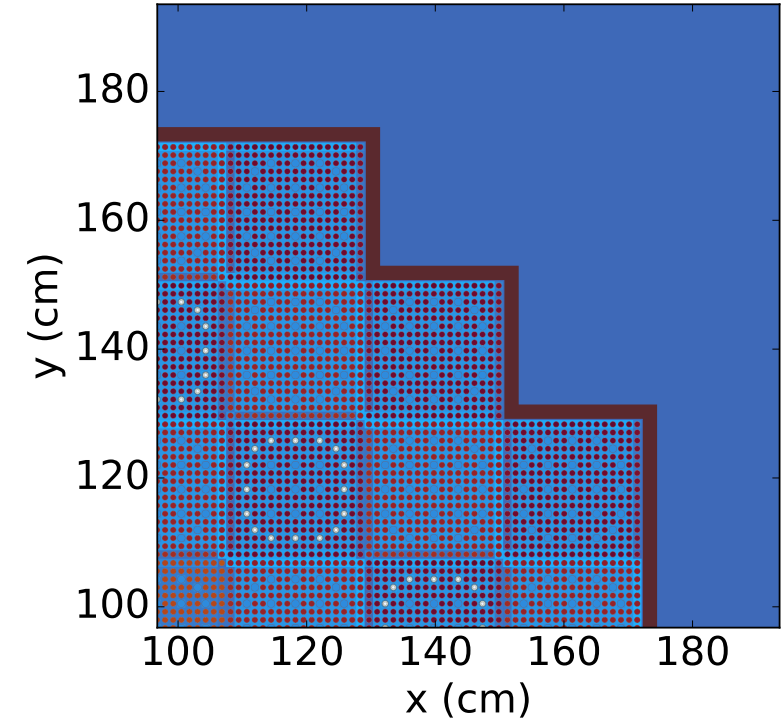
- A polynomial approximation can be used to reduce the number of $W(\cdot)$ evaluations

$$\sigma(E, T) = \frac{1}{E} \sqrt{\frac{A\pi}{k_B T} \sum_j \Re \left[r_j W \left((\sqrt{E} - p_j) \sqrt{A/k_B T} \right) \right]} + \sum_{n=0}^{N-1} a_{w,n} \mathcal{D}_n$$



GPU Performance

- Performance testing with a quarter-core model of the awaited NuScale Small Modular Reactor (SMR)
- No significant sacrifice of accuracy compared to standard continuous energy (CE) data
- Each GPU thread does individual Faddeeva evaluations (no vectorization over nuclides)
- Factor of 2-3 performance penalty on both the CPU and GPU using Pole Method for Doppler Broadening



2x IBM Power8+
4x NVIDIA Tesla P100

Questions?

