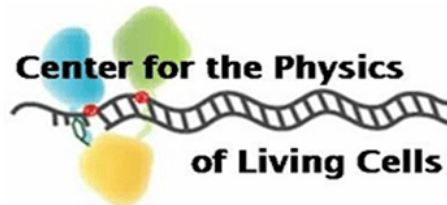


Lattice Microbes: CUDA Algorithms for Stochastic Simulation of Biochemical Reactions

Michael J. Hallock
Joseph R. Peterson
Luthey-Schulten Laboratory



GPU Programming for Molecular Modeling
Workshop, 2013

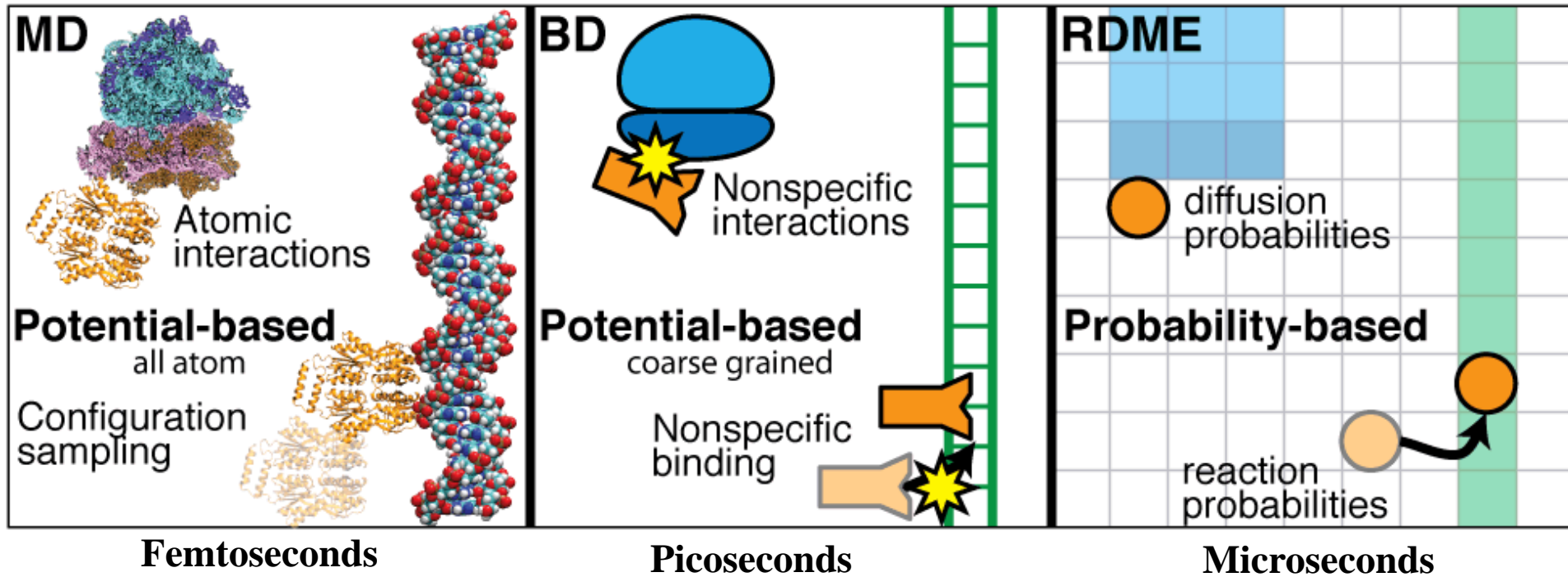
Modeling of Molecules to Cells

Biomolecular interactions span many orders of magnitude in space and time

Molecular Dynamics

Brownian Dynamics

Reactions and Diffusion

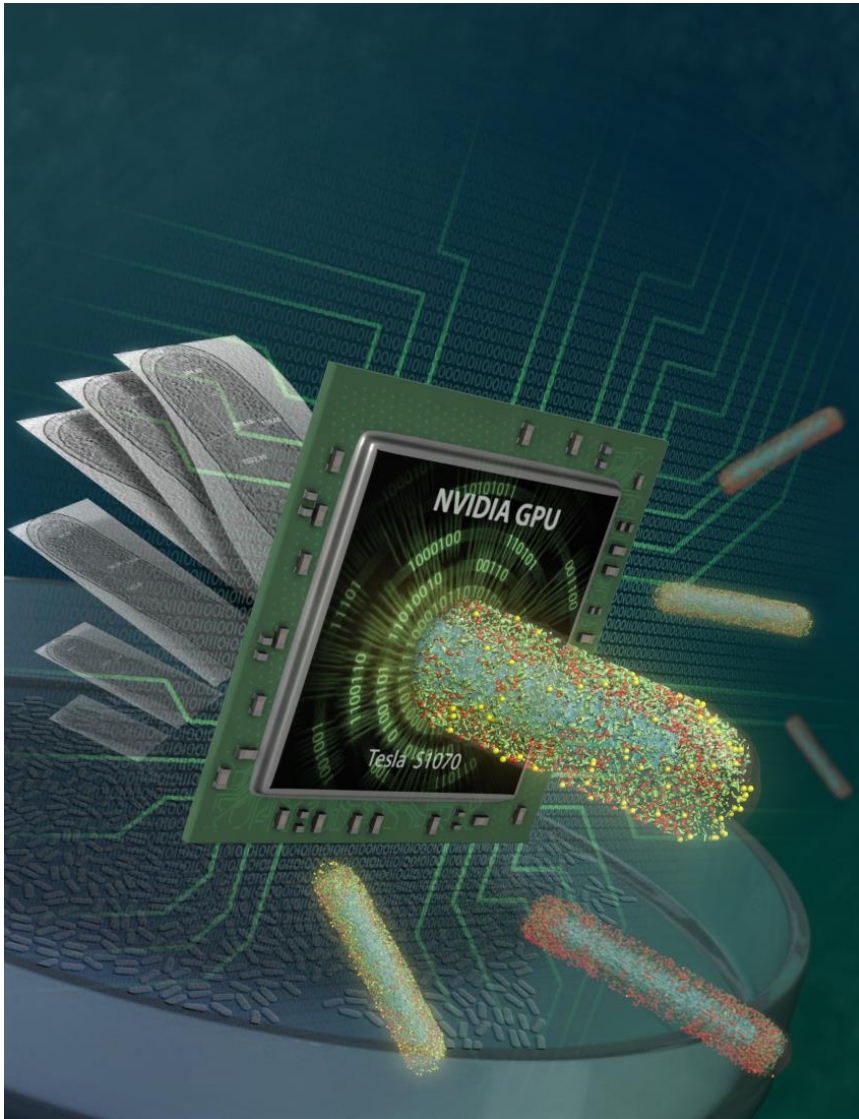


Molecules to Macromolecular assemblies

Whole Cells

Lattice Microbes

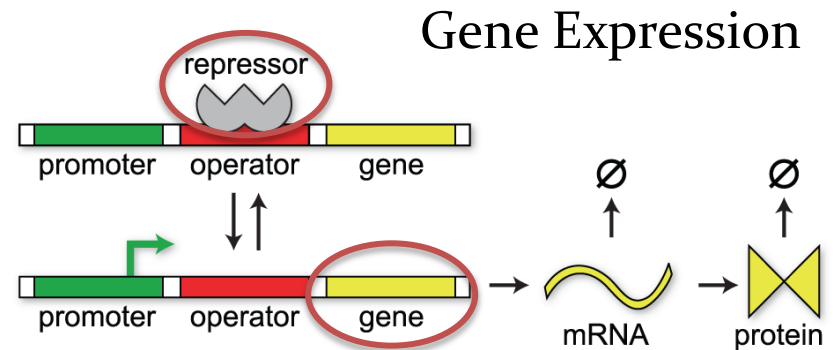
High Performance Stochastic Simulator



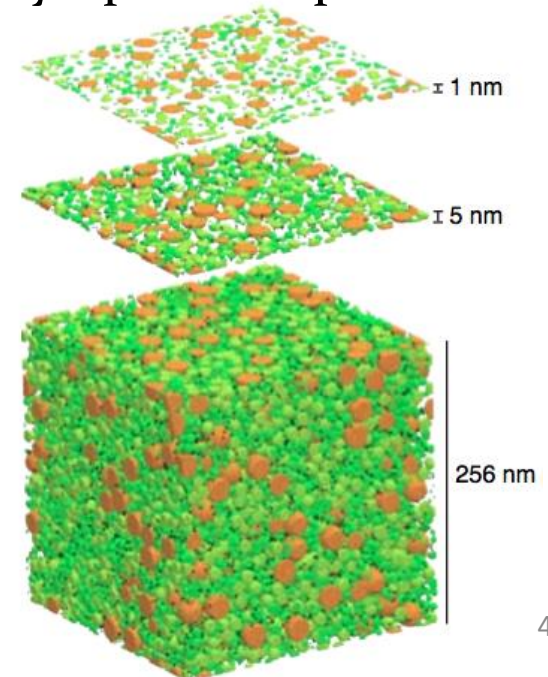
- Whole Cell Models
 - Based on experimentally determined geometry
 - Realistic *in vivo* environment
- Python interface
- Employs new algorithm to sample the Reaction-Diffusion Master Equation :
Multi-particle Diffusion (MPD-RDME)
 - Lattice based method
 - Enables acceleration by GPU hardware
 - Up to 356x faster than similar codes

Biochemical Processes Can be Stochastic

- Small copy numbers involved
 - 1 or 2 copies of gene
 - Few copies of transcription factors
- Intracellular space is very crowded
 - Not homogenous

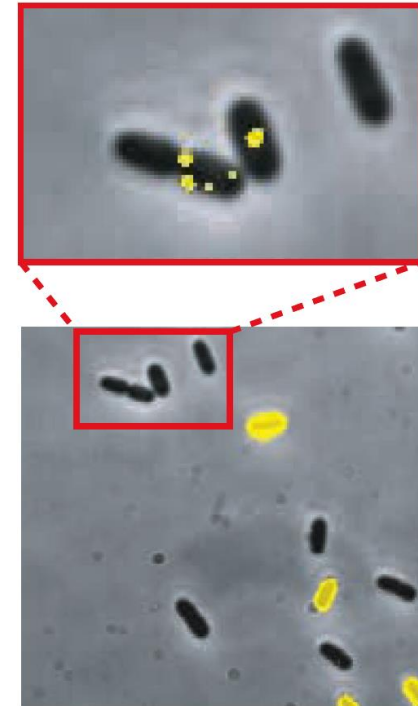
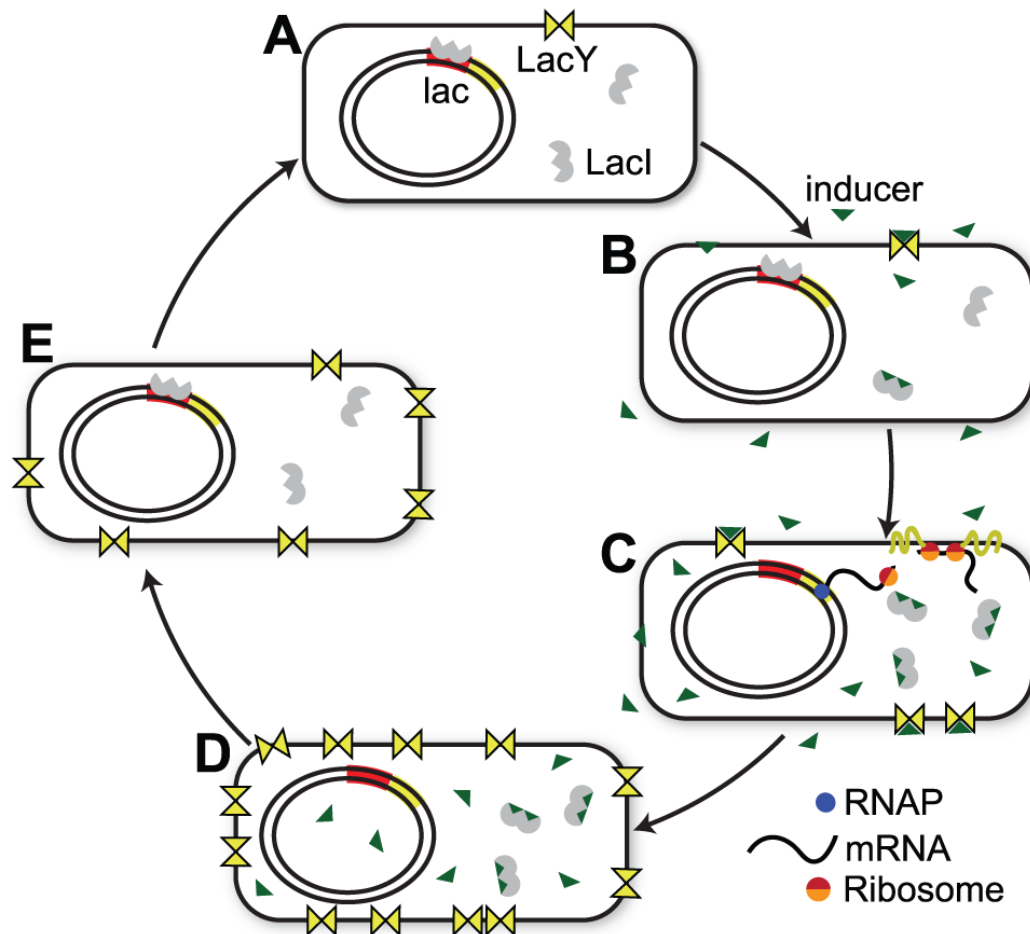


Cytoplasmic space



Stochastic Example: Two State Lac Genetic Switch

- Inducible switch for utilization of lactose by *E. coli*.



Bimodality in Population

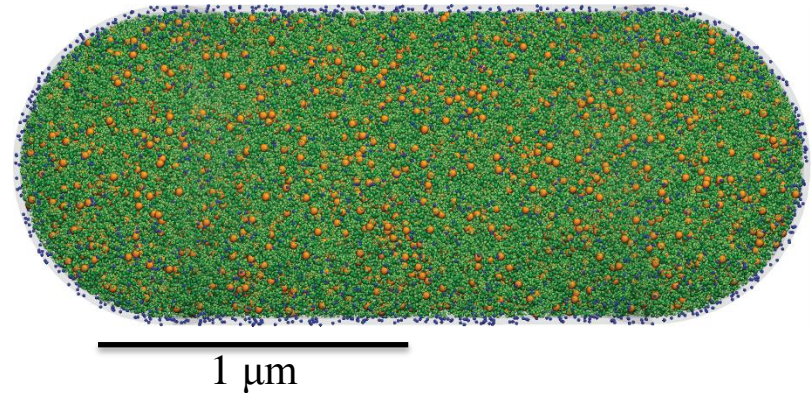
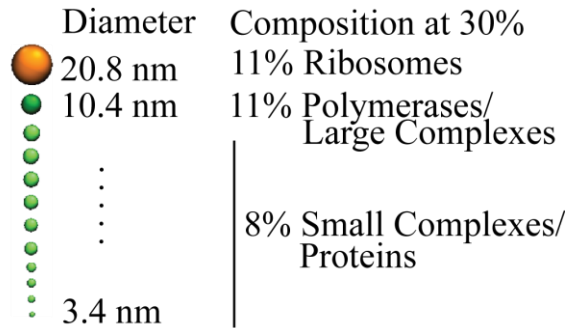
Switching due to mRNA Bursting



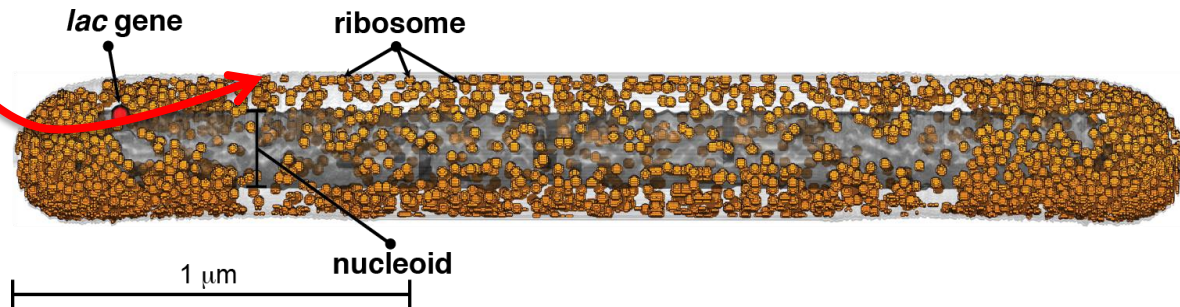
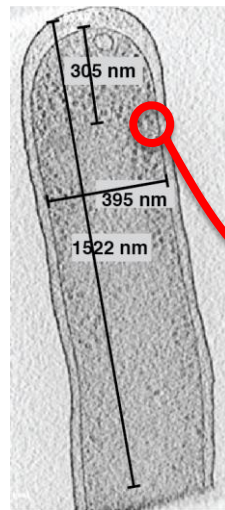
Lattice Microbes Workflow

Step 1: Making Whole Cell Models

Proteomics Data

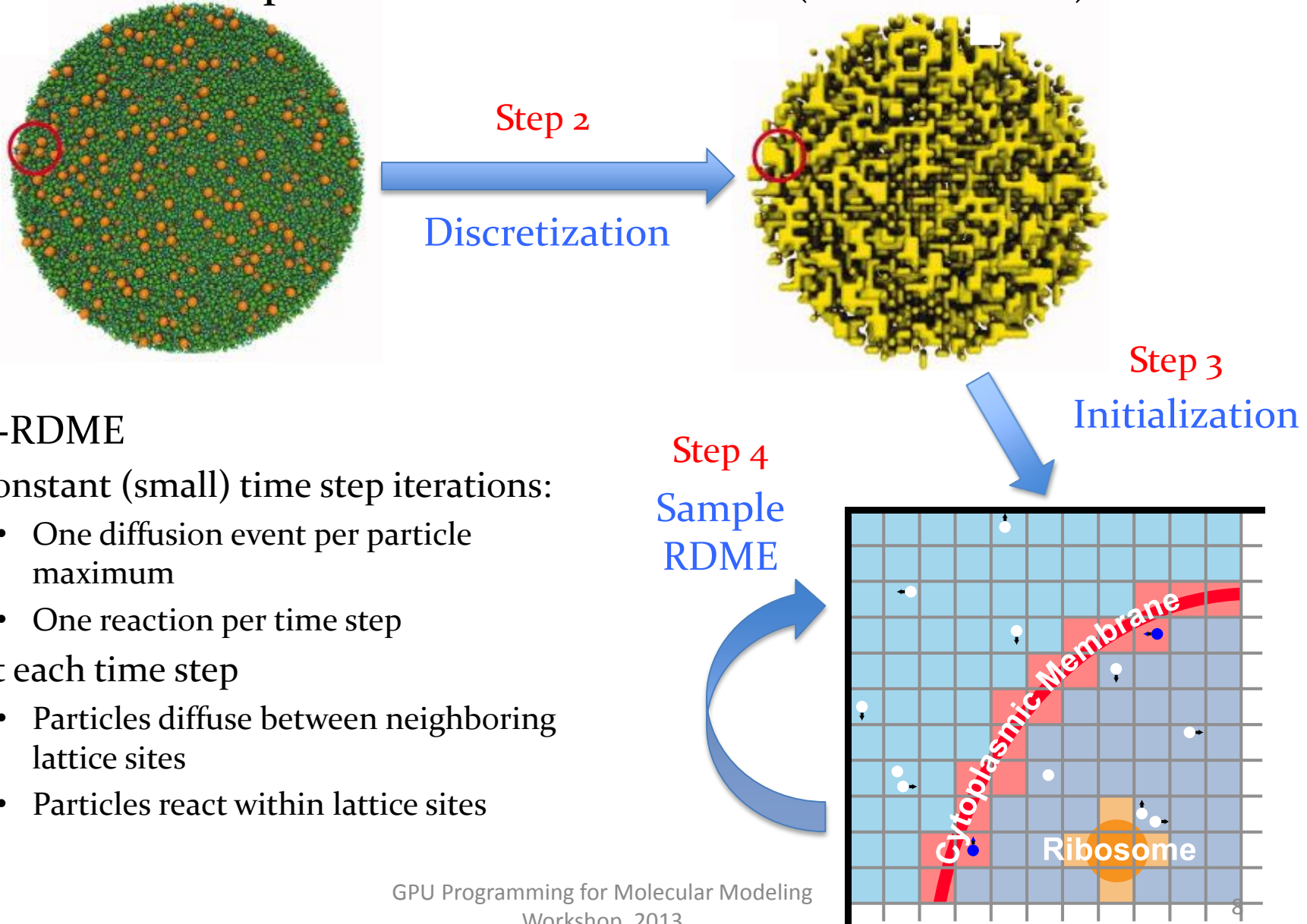


Cryo-electron Tomography



Lattice Microbes Workflow

Multi-particle Diffusion RDME (MPD-RDME)



MPD-RDME

- Constant (small) time step iterations:
 - One diffusion event per particle maximum
 - One reaction per time step
- At each time step
 - Particles diffuse between neighboring lattice sites
 - Particles react within lattice sites

Modeling Biochemical Reaction Networks

Deterministic		Stochastic (Probabilistic)	
Flux Balance Analysis (FBA)	Ordinary Differential Equations (ODE)	Chemical Master Equation (CME)	Reaction-Diffusion Master Equation (RDME)
$S \times v = 0$	$S \times v = \frac{d [C]}{dt}$	Lattice Microbes $\frac{dP(x,t)}{dt} = RP(x,t)$	$\frac{dP(x,t)}{dt} = RP(x,t) + DP(x,t)$

Full RDME:

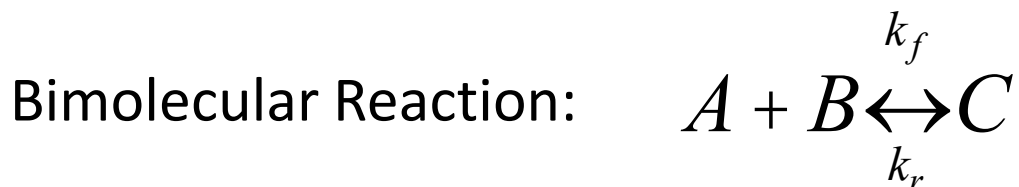
$$\frac{\partial P(x,t)}{\partial t} = \sum_{v \in V} \sum_{r=1}^R \dot{a}_r \dot{a}_r - a_r(x_v) \cdot P(x_v, t) + a_r(x_v - S_r) \cdot P(x_v - S_r, t) + \sum_{v \in V} \sum_{x \in X} \sum_{a \in A} \dot{a}_v \dot{a}_x \dot{a}_a - d^a x_v^a \cdot P(x, t) + d^a (x_{v+x}^a + 1) \cdot P(x + 1_{v+x}^a - 1_v^a, t)$$

Chemical Master Equation

$$\frac{\partial P(x, t)}{\partial t} = \sum_{r=1}^R -a_r(x)P(x, t) + a_r(x - S_r)P(x - S_r, t)$$

Reverse reaction

Forward reaction



- Probability based reaction equation
 - Applicable to discrete/stochastic (Markovian) scenarios
 - a_r linked directly to deterministic rate constant
- Applicable in well-stirred systems

$$a_r(x) = k_r x_a$$

$$a_r(x - S_r) = \frac{k_f x_a x_b}{N_A V}$$

Reaction-Diffusion Master Equation

$$\frac{\partial P(x, t)}{\partial t} = \sum_{v \in V} \sum_{r=1}^R -a_r(x_v) P(x_v, t) + a_r(x_v - S_r) P(x_v - S_r, t) + \sum_{v, x, a} \sum_{\pm i, \hat{j}, \hat{k}}^V -d^a x_v^a P(x, t) + d^a (x_{v+x}^a + 1) P(x + 1_{v+x}^a - 1_v^a, t)$$

Diffusion into subvolume

Diffusion out of subvolume

Diffusive propensity: $d^a = \frac{D^a}{l^2}$

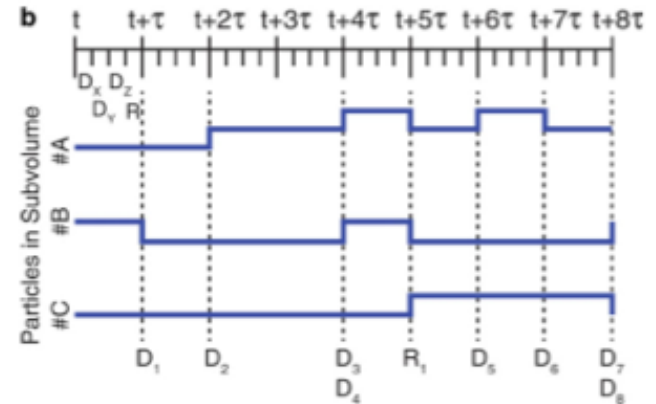
- MPD-RDME
 - Cubic subvolumes of dimension: λ
 - Each subvolume assumed to be well-stirred
 - Diffusion operations to and from nearest neighbors
 - Fine-grained parallelism

MPD-RDME Algorithm

Algorithm 1: the MPD-RDME algorithm

```

1 while  $t < t_{end}$  do
2   # Diffusion operator.
3   for  $dim \in \{x, y, z\}$  do
4     parfor  $\nu \in V$  do
5       for  $particle \in \mathbf{x}_\nu$  do
6          $n = \text{uniformRand}()$ 
7         if  $n \leq P(\text{particle}, D^{+dim})$  then
8           move  $particle$  from subvolume  $\mathbf{x}_\nu$  to neighboring subvolume in the
              + $dim$  direction
9         else if  $(n - P(\text{particle}, D^{+dim})) \leq P(\text{particle}, D^{-dim})$  then
10          move  $particle$  from subvolume  $\mathbf{x}_\nu$  to neighboring subvolume in the
              - $dim$  direction
11        end
12      end
13    end
14  end
15  # Reaction operator.
16  parfor  $\nu \in V$  do
17    for  $r \in R$  do
18       $a_{tot} = a_{tot} + a_r(\mathbf{x}_\nu)$ 
19    end
20     $n_1 = \text{uniformRand}()$ 
21    if  $n_1 \leq \int_0^\tau a_{tot} e^{-a_{tot}t} dt$  then
22       $n_2 = \text{uniformRand}()$ 
23      for  $r \in R$  do
24        if  $a_{r-1}(\mathbf{x}_\nu) < n_2 \cdot a_{tot} \leq a_r(\mathbf{x}_\nu)$  then
25          perform reaction  $r$  in subvolume  $\mathbf{x}_\nu$ 
26        end
27      end
28    end
29  end
30   $t = t + \tau$ 
31 end
  
```

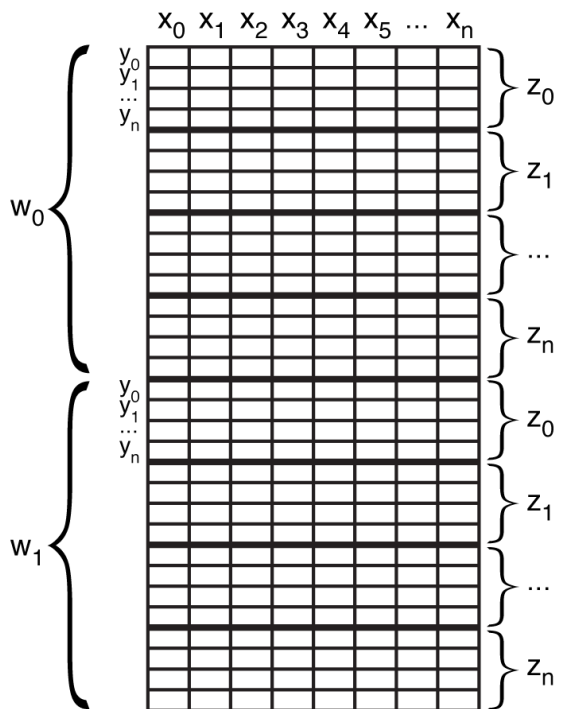
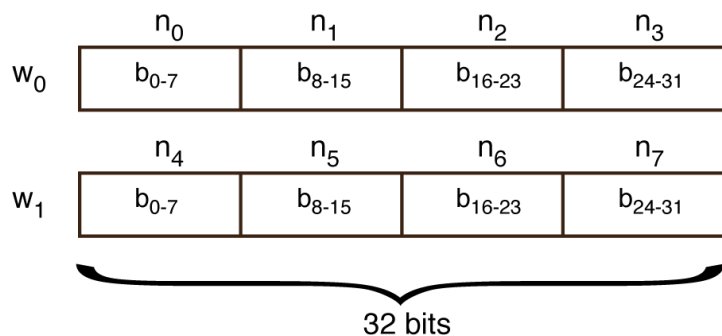


- Many short time steps:

$$t \in \frac{l^2}{6D_{\max}}$$

- GPU Implementation
 - Three sequential diffusion kernels: x, y, z
 - One reaction kernel

Lattice Data Structure



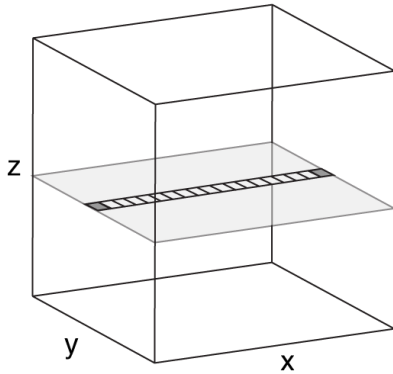
$$\mathcal{L}_{\text{index}}^{0-3} = x + (y \cdot x_n) + (z \cdot x_n \cdot y_n)$$

$$\mathcal{L}_{\text{index}}^{4-7} = x + (y \cdot x_n) + (z \cdot x_n \cdot y_n) + (x_n \cdot y_n \cdot z_n)$$

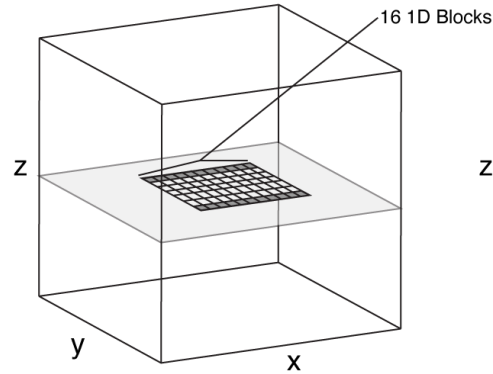
- 2 Lattices mapped as a 3D array in linear memory
 - Site type lattice
 - Particle lattice
 - 8 bit particles \rightarrow 256 particle types
 - 8 particles per site \rightarrow 320 μM at 10% full for 16 nm spacing

Diffusion Kernels

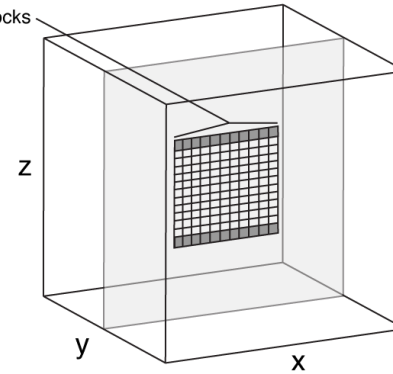
Thread Block Pattern
(X Diffusion)



Thread Block Pattern
(Y Diffusion)



Thread Block Pattern
(Z Diffusion)



Coalesced memory access:

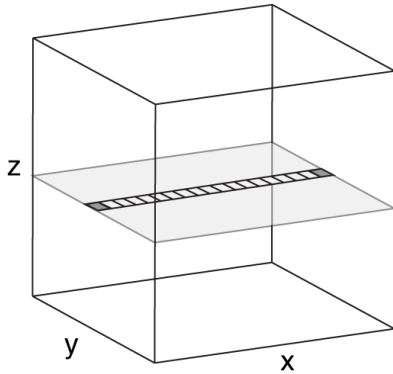
- 1D x axis
- 2D x-y and x-z blocks

Diffusion Operation

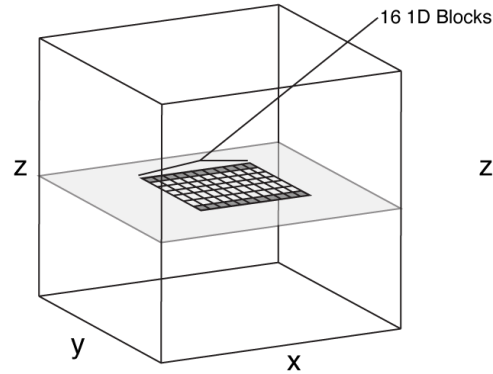
```
__global__ void MpdRdmeSolver_x_kernel(...)  
{  
    // Figure out the offset of this thread and block in the lattice and the lattice segment.  
    ...  
    // Shared memory to store the lattice segment.  
    __shared__ unsigned int window[MPD_X_WINDOW_SIZE*MPD_WORDS_PER_SITE];  
    __shared__ uint8_t sitesWindow[MPD_X_WINDOW_SIZE];  
  
    // Copy the x window from device memory into shared memory.  
    copyXWindowFromLattice(bx, inLattice, window, latticeIndex, latticeXIndex, windowIndex);  
    copyXWindowFromSites(bx, inSites, sitesWindow, latticeIndex, latticeXIndex, windowIndex);  
    __syncthreads();  
  
    // Make the choices.  
    __shared__ unsigned int choices[MPD_X_WINDOW_SIZE*MPD_WORDS_PER_SITE];  
    makeXDiffusionChoices(window, sitesWindow, choices...);  
    __syncthreads();  
  
    // Propagate the choices to the new lattice segment.  
    performPropagation(...);  
}
```

Diffusion Kernels

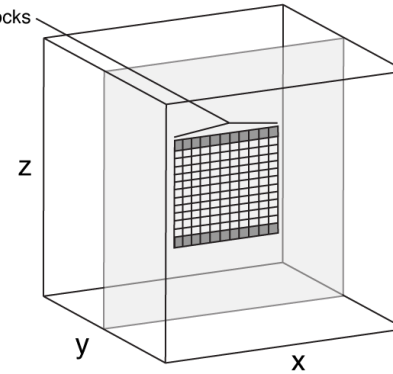
Thread Block Pattern
(X Diffusion)



Thread Block Pattern
(Y Diffusion)



Thread Block Pattern
(Z Diffusion)



Coalesced memory access:

- 1D x axis
- 2D x-y and x-z blocks

X-diffusion copy lattice window:

```
inline __device__ void copyXWindowFromLattice(...)
{
    if (latticeXIndex < latticeXSizeC)
    {
        // Load the block.
        for (uint w=0, latticeOffset=0, windowOffset=0;
             w < MPD_WORDS_PER_SITE;
             w++, latticeOffset+=latticeXYZSizeC, windowOffset+=MPD_X_WINDOW_SIZE)
            window[windowIndex+windowOffset] = lattice[latticeIndex+latticeOffset];

        // Calculate the effective thread block width,
        // accounting for the fact that the block might not be aligned with the end of the lattice.
        int threadBlockWidth = ((bx+1)*blockDim.x <= latticeXSizeC)?(blockDim.x):(latticeXSizeC-(bx*blockDim.x));

        // Load leading apron
        ...
        // Load trailing apron
        ...
    }
}
```

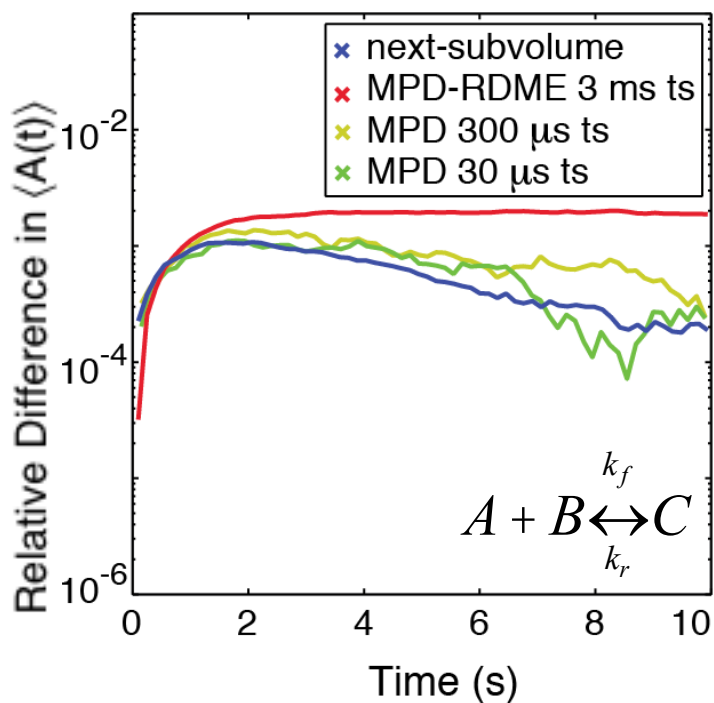
Code Performance

Single GPU version 2.0 available since January 2012

Suitable error from MPD-RDME approximation:

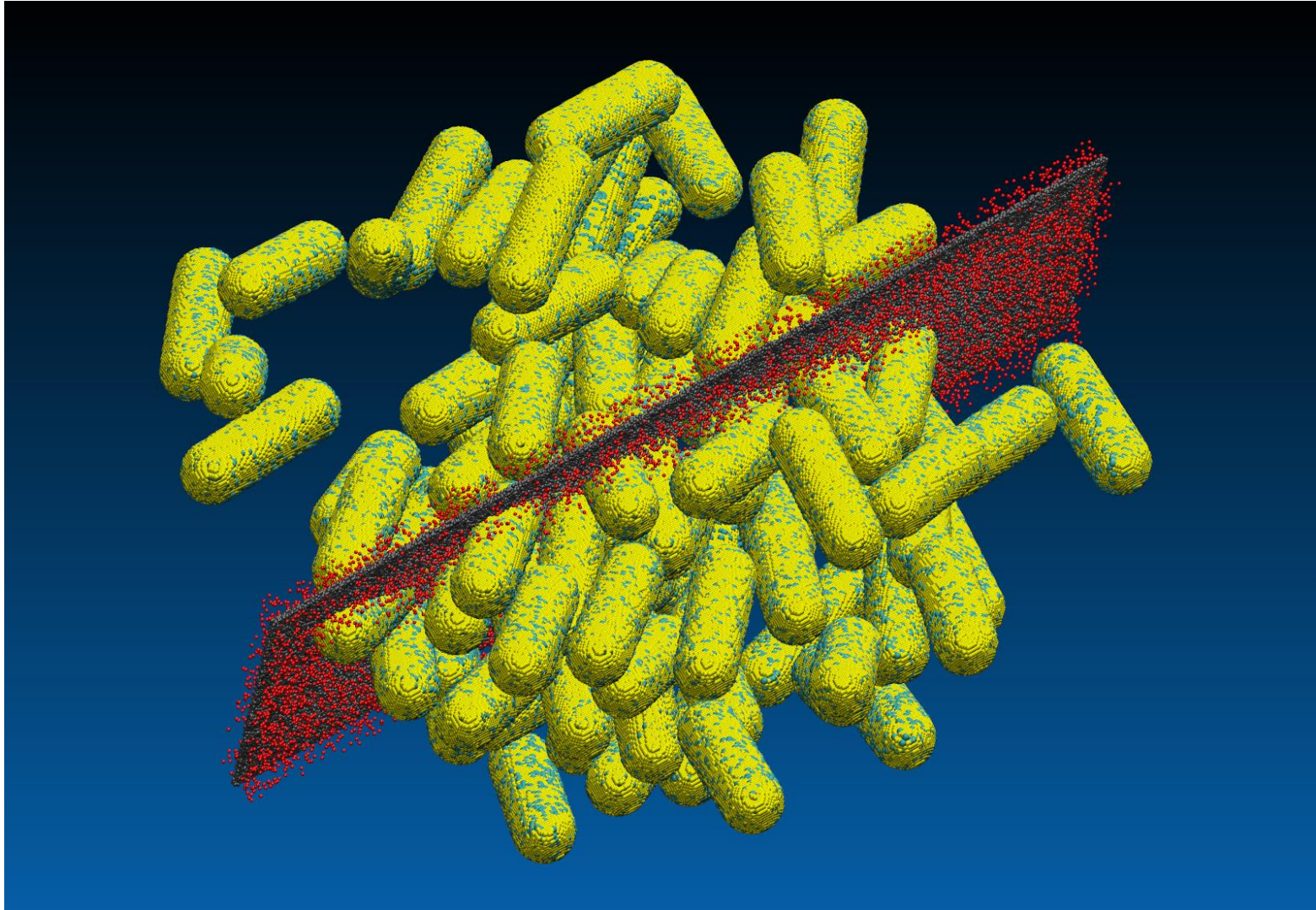
RDME: as much as *356 times faster*

CME: as much as *2.7 times faster*



Method	Lac Switch, [I] = 5 μ M		Lac Switch, [I] = 40 μ M	
	rxns/sec	wall clock/sim hr	rxns/sec	wall clock/sim hr
CME				
Direct GPU	7.3x10 ⁶	0.018 sec	7.4x10 ⁶	4.0 sec
Direct No-GPU	5.0x10 ⁶	0.026 sec	5.0x10 ⁶	5.9 sec
Next React GPU	5.5x10 ⁶	0.024 sec	6.2x10 ⁶	4.8 sec
Next React No-GPU	3.9x10 ⁶	0.032 sec	4.5x10 ⁶	6.5 sec
Stochkit SSA	4.2x10 ⁶	0.031 sec	4.2x10 ⁶	7.2 sec
COPASI	2.7x10 ⁶	0.048 sec	2.8x10 ⁶	11 sec
RDME				
MPD-RDME GPU	0.44	1.89 days	0.1200	0.191 yrs
Next Subvol GPU	0.10	8.33 days	0.0004	59 yrs
Next Subvol No-GPU	0.10	8.67 days	0.0004	59 yrs
MesoRD	0.01	83 days	0.0003	68 yrs

Larger and More Complex Models: Motivating Multiple GPUs



Multiple GPU Communication

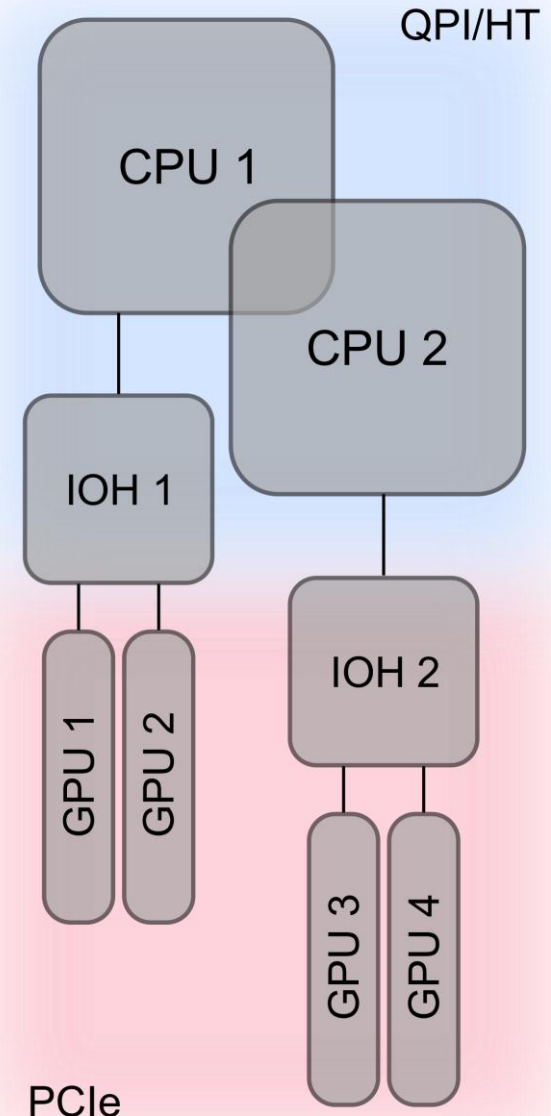
Data from one GPU to another can be moved:

- From GPU 1 to host memory, and then from host memory to GPU 2
- DMA directly from GPU 1 to GPU 2

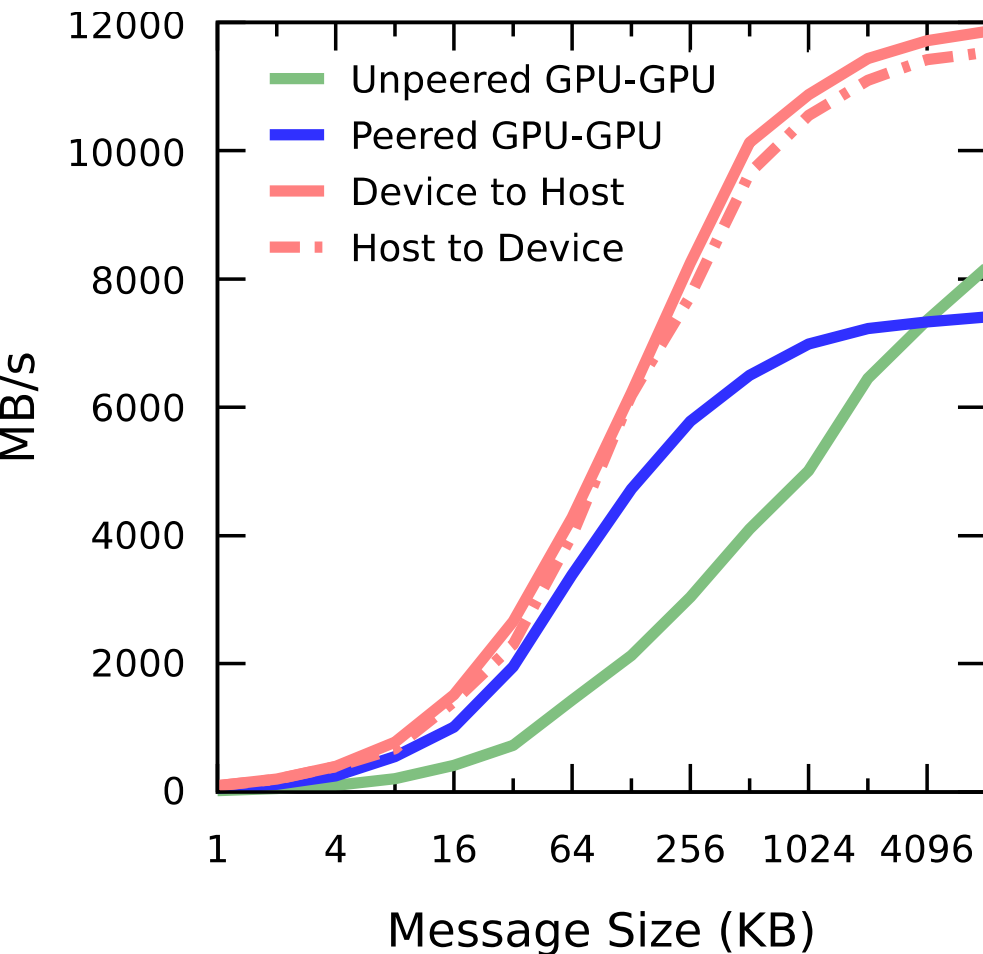
Device-to-Device copies are accelerated if the GPUs are *peered* and can directly address each other's memory

Ability to peer devices depends on machine topology, multi-IOH machines can result in GPU in non-peering GPUs

Check using `cudaDeviceCanAccessPeer()`



Bandwidth of Peered GPUs



Peer copies alleviate pressure on host memory by only occupying PCIe bus

Can allow for greater aggregate bandwidth

Test shown using two GTX680 cards on a Sandy Bridge system running at PCIe 3.0 speeds, which is not officially supported.

Peering GPUs

Peering is one-way, so two calls to `cudaDeviceEnablePeerAccess` are required to have bi-directional communication between devices

Thread 1:

```
cudaSetDevice(1);  
cudaDeviceEnablePeerAccess(2, 0); // 0 = flags (reserved)
```

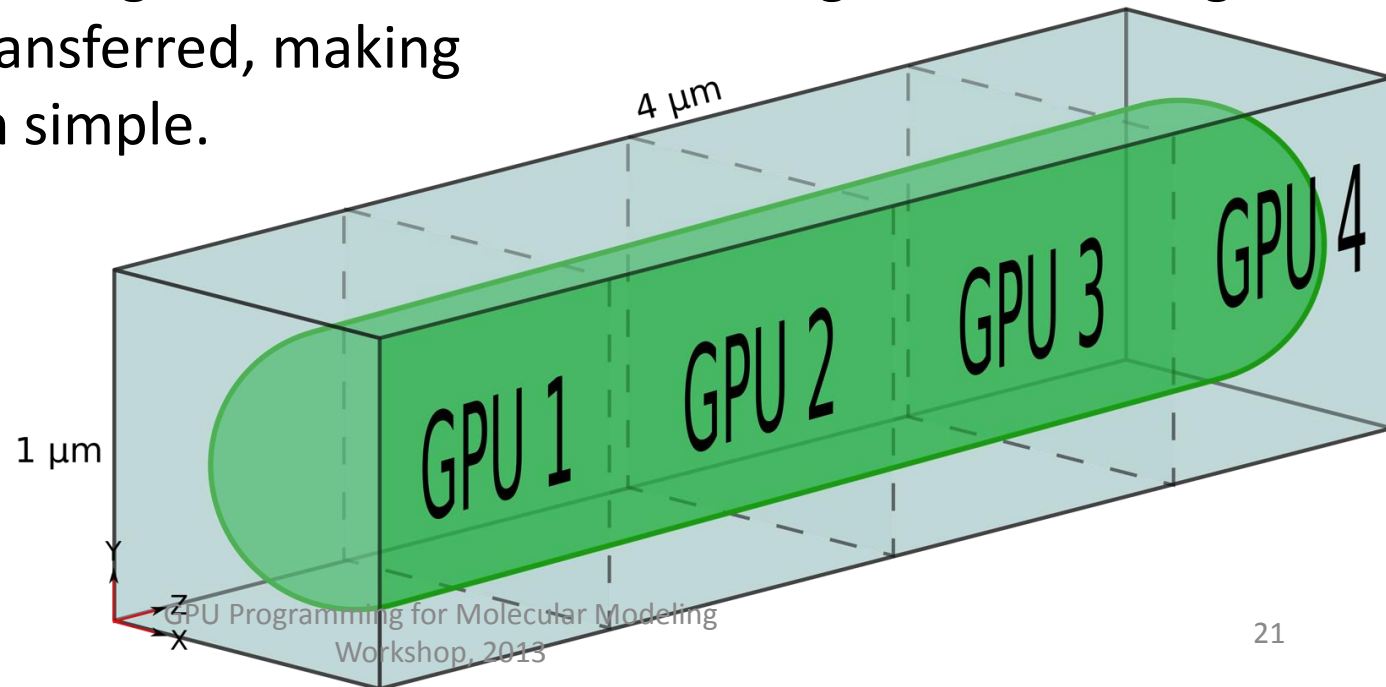
Thread 2:

```
cudaSetDevice(2);  
cudaDeviceEnablePeerAccess(1, 0);
```

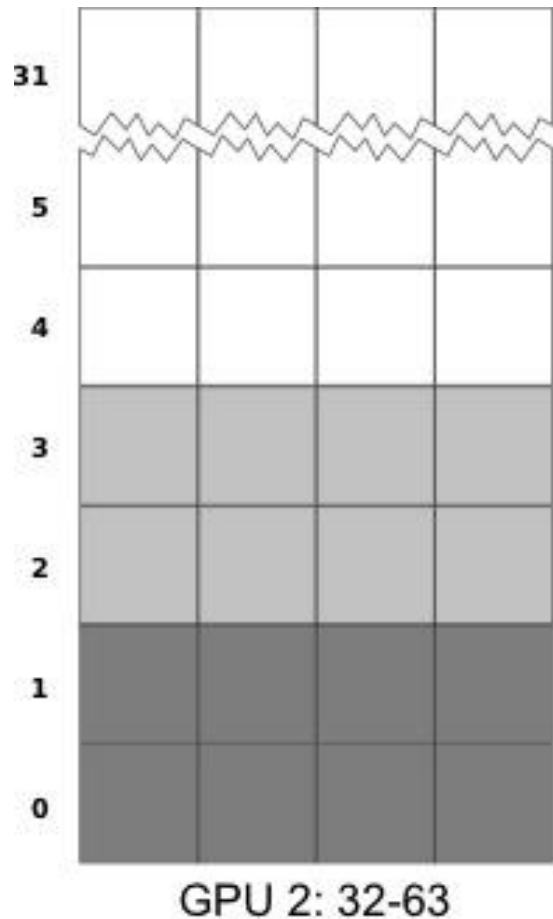

MPD-RDME over Multiple GPUs

To use multiple GPUs, we must divide up the work.

Apply a spatial decomposition to slice the lattice into sub-lattices to assign to each available GPU. A one-dimensional decomposition along the z-axis allows for contiguous data ranges that must be transferred, making communication simple.

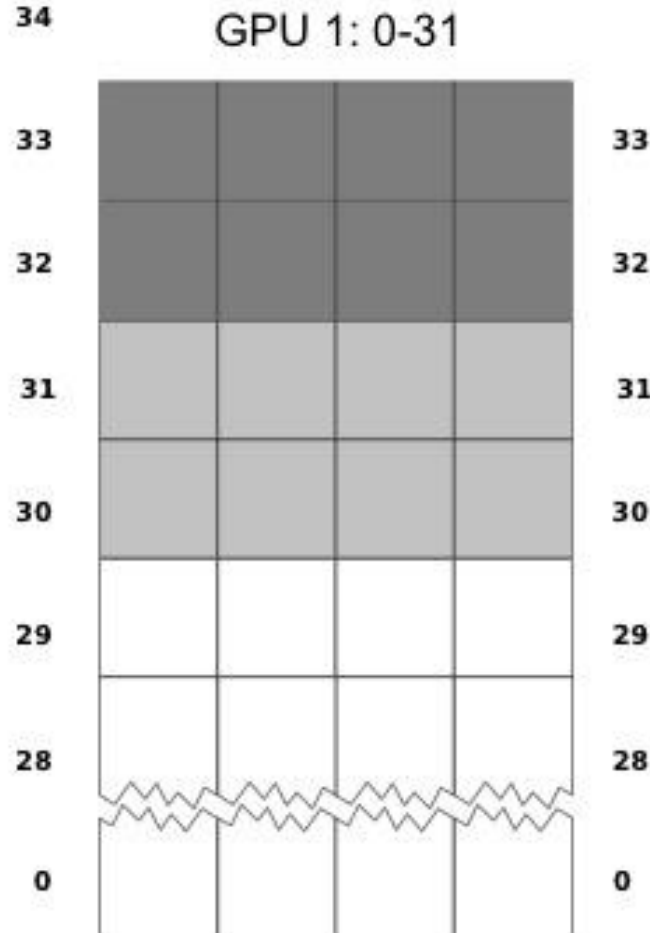


Inter-GPU Data Dependencies



63
35
34
GPU 1: 0-31

Illustration of a GPU boundary for an X×Y×64 lattice

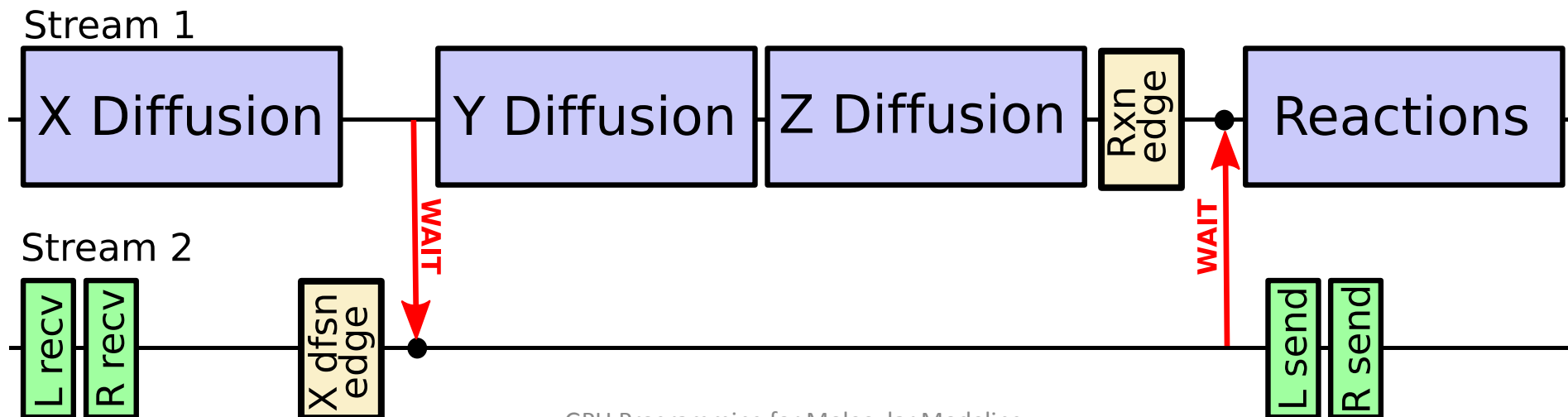


Dark gray regions are the extensions to cover the neighboring lattice, and must be updated from the neighbor every timestep.

Light gray regions are data that needs a neighbor depends upon.

Hiding Copy Latency

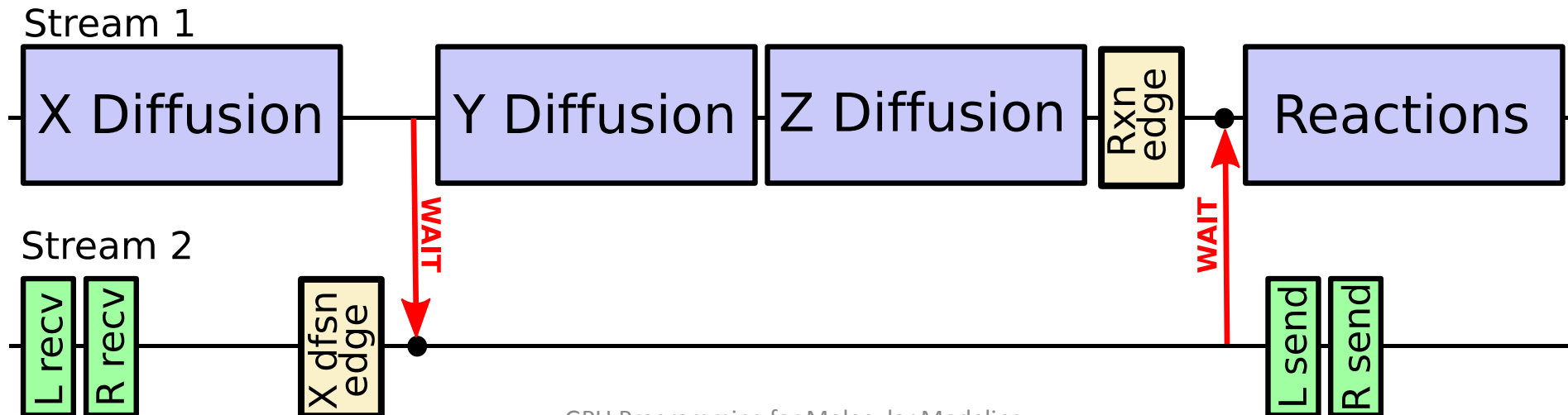
- Overlap non-dependent computation with communication using multiple CUDA streams
- Synchronize streams using events



Hiding Copy Latency

```
x_kernel<<<inner_grid,threads,0,stream1>>>( ... );  
cudaMemcpyAsync( ... , stream2 );  
x_kernel_edge<<<edge_grid, threads, 0, stream2>>>( ... );  
cudaEventRecord(x_event, stream2);  
cudaStreamWaitEvent(stream1, x_event, 0);
```

```
y_kernel<<<grid, threads, 0, stream1>>>( ... );
```

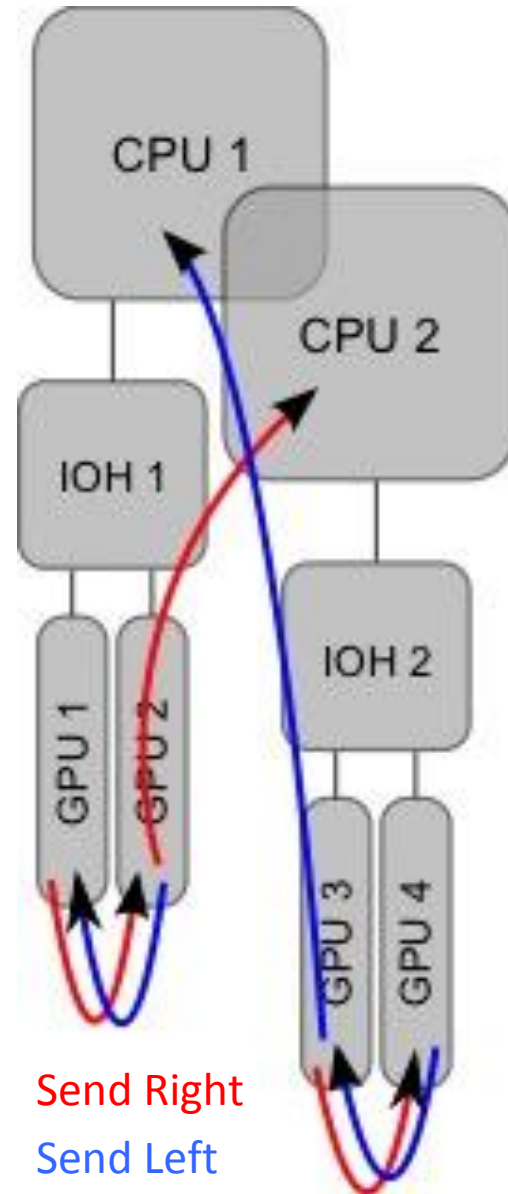


Decoupled Send and Receive

Data sent is not used until the beginning of the next timestep, so we do not need to place the updated lattice information onto the neighboring GPU if not peered.

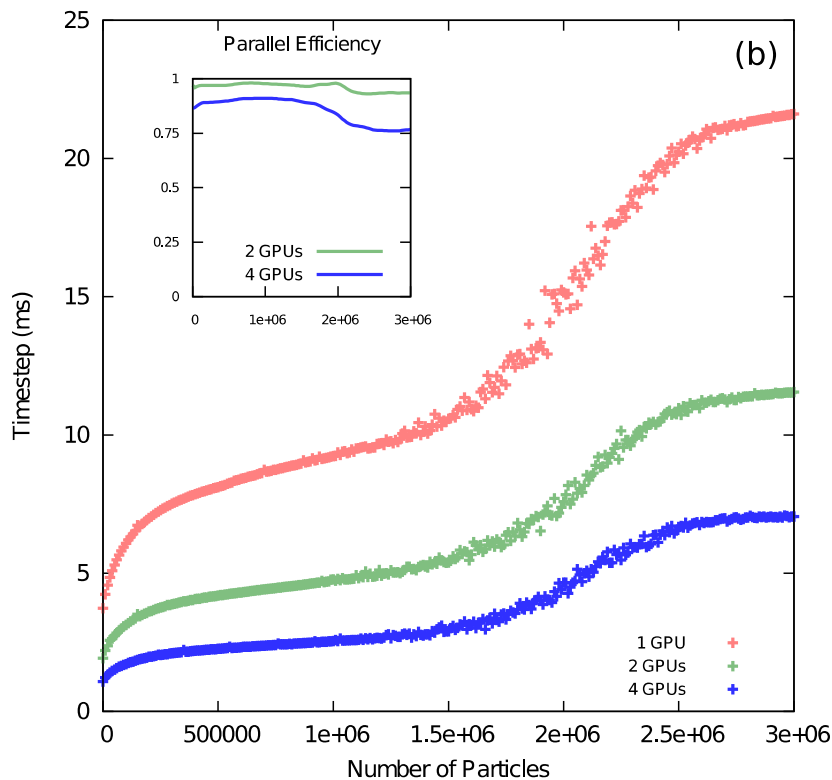
The receive operation for peered GPUs is a local memory copy, unpeered GPUs read from (NUMA local) host memory.

UVA makes it easy to implement both copies by only varying where the destination pointer was allocated.

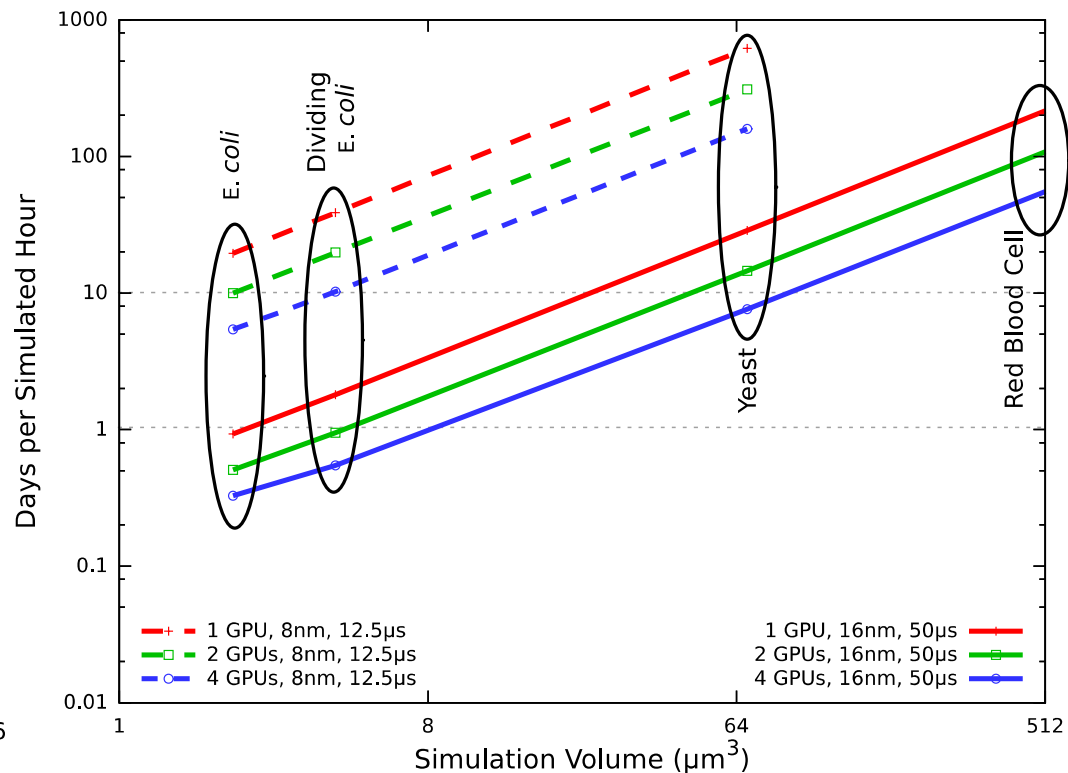


Multi-GPU Performance

Timestep and Efficiency for 128x128x256 Lattice



Benchmark System Runtimes - Eir (Four GTX680s)



Acknowledgements

P. Labhsetwar



Elijah Roberts



Zan Schulten

John



John Cole

