

Scalable AnyDSL Molecular Dynamics application with MPI

Scaling the Molecular Dynamics implementation in AnyDSL

Rafael Ravedutti Lucio Machado

Chair of Computer Science 10, Friedrich-Alexander University of Erlangen-Nuremberg

March 19, 2019



Motivation

Compare AnyDSL scalable application with other state-of-the-art technologies in order to explore benefits in code writing and performance.

Outline

AnyDSL

Molecular Dynamics

Proposal

Experimental Results

AnyDSL



AnyDSL

Framework for development of domain-specific libraries

- Higher-order functions
- Thorin
- Impala
- Partial evaluation

AnyDSL

```
1  fn main() {  
2      let img = load("dragon.png");  
3      let blurred = gaussian_blur(img);  
4  }
```

AnyDSL

```
1  fn gaussian_blur(field: Field) -> Field {
2      let stencil: Stencil = { /* ... */ };
3      let mut out: Field    = { /* ... */ };
4
5      for x, y in @iterate(out) {
6          out.data(x, y) = apply_stencil(x, y, field, stencil);
7      }
8
9      out
10 }
```

AnyDSL

```
1  fn iterate(field: Field, body: fn(int, int) -> ()) -> () {  
2      let grid = (field.cols, field.rows, 1);  
3      let block = (128, 1, 1);  
4  
5      with nvvm(grid, block) {  
6          let x = nvvm_tid_x() + nvvm_ntid_x() * nvvm_ctaid_x();  
7          let y = nvvm_tid_y() + nvvm_ntid_y() * nvvm_ctaid_y();  
8          body(x, y);  
9      }  
10 }
```


Molecular Dynamics



Molecular Dynamics

Pair-wise interaction of particles simulation implemented in AnyDSL

- Cells of particles (bounding boxes)
- Neighborlists
- Cluster of particles
- Target CPU with vectorization instructions and GPU

Molecular Dynamics

Steps

1. Initialize grid
2. Initialize clusters
3. Build neighbor lists
4. Compute forces and update particles (for 20 timesteps)
5. Redistribute particles and go back to item 2

Molecular Dynamics

Steps

1. Initialize grid
2. Initialize clusters
3. Build neighbor lists
4. Compute forces and update particles (for 20 timesteps)
5. Redistribute particles and go back to item 2

Molecular Dynamics

Steps

1. Initialize grid
2. Initialize clusters
3. Build neighbor lists
4. Compute forces and update particles (for 20 timesteps)
5. Redistribute particles and go back to item 2

Molecular Dynamics

Steps

1. Initialize grid
2. Initialize clusters
3. Build neighbor lists
4. Compute forces and update particles (for 20 timesteps)
5. Redistribute particles and go back to item 2

Molecular Dynamics

Steps

1. Initialize grid
2. Initialize clusters
3. Build neighbor lists
4. Compute forces and update particles (for 20 timesteps)
5. Redistribute particles and go back to item 2

Proposal



Proposal

Goals

- Scalable version of the application
- First: scale application on homogeneous clusters
- In the future: heterogeneous clusters (both CPU and GPU nodes)
- Compare scalable implementation with other state-of-the-art versions

Proposal

Steps

- Domain partitioning
- Communication pattern
- Synchronization of cells (every timestep)
- Particle exchange (after redistribution)

Proposal

Domain partitioning

- Define configuration of nodes
- Split domain accordingly
- Define current node bounding box
- Include ghost layer cells
- **get_sync_timesteps()**

Domain partitioning

```

1  fn get_node_config(
2      world_size: i32,
3      rank: i32,
4      xcells: i32,
5      ycells: i32,
6      zcells: i32) -> [i32 * 3] {
7
8      let mut gx = 1, gy = 1, gz = 1;
9      let mut min_missing_factor = xcells * ycells * zcells;
10
11     for i in range(1, world_size) {
12         if(world_size % i == 0) {
13             let rem_yz = world_size / i;
14
15             for j in range(1, rem_yz) {
16                 if(rem_yz % j == 0) {
17                     let k = rem_yz / j;
18                     let missing_factor = xcells % i + ycells % j + zcells % k;
19
20                     if(min_missing_factor > missing_factor) {
21                         gx = i;
22                         gy = j;
23                         gz = k;
24                         min_missing_factor = missing_factor;
25                     }
26                 }
27             }
28         }
29     }
30
31     [gx, gy, gz]
32 }

```

Domain partitioning

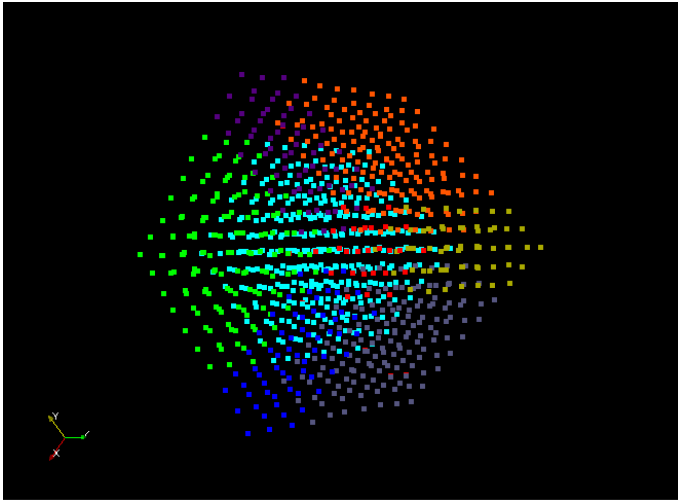
```

1  fn @get_node_bounding_box(
2      world_size: i32,
3      rank: i32,
4      cell_spacing: f64,
5      aabb: AABB) -> AABB {
6
7      let mut xmin, xmax, ... : f64;
8
9      let xcells = math.floor((aabb.max(0) - aabb.min(0)) / cell_spacing);
10     let xlength = (xcells / node_dims(0)) * cell_spacing;
11     let rank_index = unflat_index(rank, ...);
12
13     xmin = aabb.min(0) + xlength * rank_index(0);
14     xmax = aabb.min(0) + xlength * (rank_index(0) + 1);
15
16     if(rank_index(0) > 0) {
17         xmin -= get_sync_timesteps() * cell_spacing;
18     }
19
20     if(rank_index(0) < node_dims(0) - 1) {
21         xmax += get_sync_timesteps() * cell_spacing;
22     }
23
24     /* y and z are analogous to x */
25
26     AABB {
27         min: [xmin, ymin, zmin],
28         max: [xmax, ymax, zmax]
29     }
30 }

```

Proposal

Domain partitioning



Proposal

Communication pattern

- Higher-order function for iteration
- Easy to write and change with AnyDSL

Communication pattern

```
1  fn communication_nodes (  
2      world_size: i32,  
3      rank: i32,  
4      grid: Grid,  
5      body: fn(i32, [i32 * 3], [i32 * 3], [i32 * 3], [i32 * 3]) -> () -> () {  
6  
7      ...  
8  
9      if(rank_index(0) < node_dims(0)) {  
10         @body(  
11             flat_index([rank_index(0) + 1, rank_index(1), rank_index(2)], ...),  
12             send_begin1, send_end1, recv_begin1, recv_end1 // regions to communicate  
13         );  
14     }  
15  
16     if(rank_index(0) > 0) {  
17         @body(  
18             flat_index([rank_index(0) - 1, rank_index(1), rank_index(2)], ...),  
19             send_begin2, send_end2, recv_begin2, recv_end2 // regions to communicate  
20         );  
21     }  
22  
23     ... /* y and z are analogous to x */  
24 }
```


Proposal

Synchronization of cells

- Update positions, velocity and forces of particles
- One-step communication
- Every **get_sync_timesteps()** timesteps

Synchronization of cells

```
1  fn synchronize_ghost_layer_cells(  
2      grid: &mut Grid,  
3      accelerator_grid: AcceleratorGrid,  
4      world_size: i32,  
5      world_rank: i32) -> () {  
6  
7      ... /* Transfer data from accelerator to CPU */  
8  
9      for exchange_rank, send_begin, send_end, recv_begin, recv_end in  
10         communication_nodes(world_size, world_rank, *grid) {  
11  
12         pack_ghost_layer_cells(..., send_begin, send_end);  
13  
14         mpih.irecv(...);  
15         mpih.send(...);  
16         mpih.wait(...);  
17  
18         unpack_ghost_layer_cells(..., recv_begin, recv_end);  
19     }  
20  
21     ... /* Transfer data from CPU to accelerator */  
22 }
```

Proposal

Particle exchange

- Exchange redistributed cells
- May be a two-step communication ($N' > N + N/2$)
- Every 20 timesteps (redistribution)

Particle exchange

```
1  fn exchange_ghost_layer_particles(  
2      grid: &mut Grid,  
3      world_size: i32,  
4      world_rank: i32) -> () {  
5  
6      ...  
7  
8      for exchange_rank, send_begin, send_end, recv_begin, recv_end in  
9          communication_nodes(world_size, world_rank, *grid) {  
10         ...  
11  
12         pack_ghost_layer_particles(..., &mut rmng_send_ptcs);  
13         mpih.irecv(...);  
14         mpih.send(...);  
15         mpih.wait(...);  
16         unpack_ghost_layer_particles(..., &mut rmng_recv_ptcs);  
17  
18         if (rmng_recv_ptcs > 0) {  
19             mpih.irecv(...);  
20         }  
21  
22         if (rmng_send_ptcs > 0) {  
23             pack_ghost_layer_particles(...);  
24             mpih.send(...);  
25         }  
26  
27         if (rmng_recv_ptcs > 0) {  
28             mpih.wait(...);  
29             unpack_ghost_layer_particles(...);  
30         }  
31     }  
32 }
```

Experimental Results



Experimental Results

Cluster configuration

- 36-CPU High-Performance-Cluster
- 8 compute nodes
- 4 x Intel(R) Xeon(R) CPU E7-4830, 2.13 GHz 2.4GHz (max. turbo) (8 cores + SMT), SSE 4.1/4.2, 24 MB shared cache
- 256 GB RAM
- 2 x 300 GB SAS internal disks
- NVIDIA GeForce GTX 680
- QDR Infiniband network

Experimental Results

Cluster mapping

- MPI processes are mapped to different CPU cores in the same cluster node
- We bind processes to sockets for better memory usage

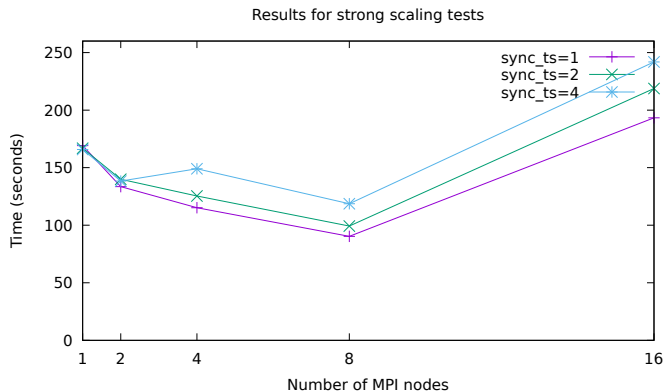
Experimental Results

Strong scaling

- Grid size: 128x128x128
- About 2 million particles

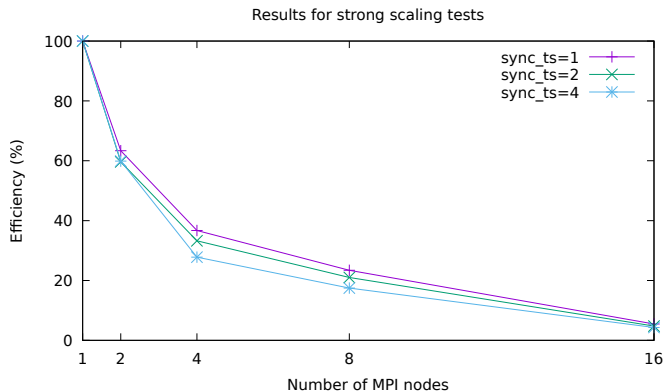
Experimental Results

Strong scaling time



Experimental Results

Strong scaling efficiency



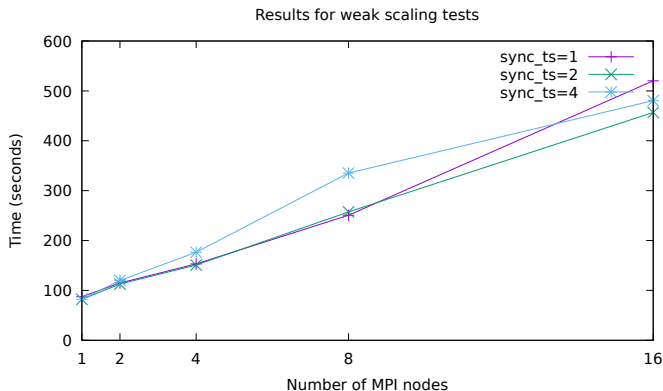
Experimental Results

Weak scaling

- **Problem: cannot use enough particles (too much memory)**
- Must fix this in order to get good weak scaling results

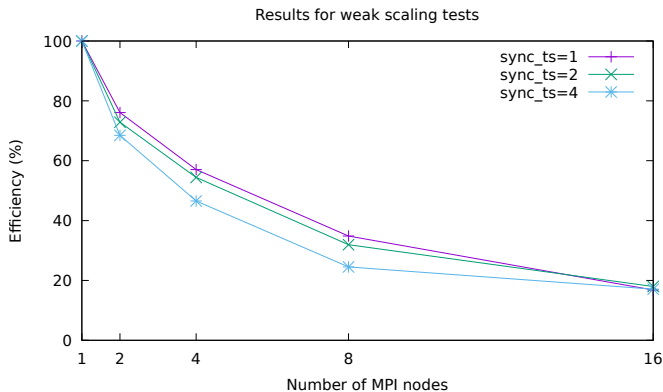
Experimental Results

Weak scaling time



Experimental Results

Weak scaling efficiency



Future work

- Reduce memory usage by application
- Improve partitioning for heterogeneous clusters
- Compare with state-of-the-art applications

Thanks for listening.
Any questions?

References



References I

- [1] J. Schmitt, H. Kostler, J. Eitzinger, and R. Membarth, “Unified code generation for the parallel computation of pairwise interactions using partial evaluation”, *17th International Symposium on Parallel and Distributed Computing*, 2018.