# OpenACC - Performance and Productivity

Paul Springer

Aachen Institute for Advanced Study in
Computational Engineering Science
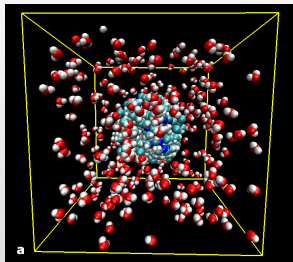
Aachen, 14.03.13 – PPCES 2013

**RWTH**AACHEN
UNIVERSITY

# Molecular Dynamics Simulation

- System of $N$ interacting particles
  - E.g.: Atoms, molecules, planets

- Simulate their motion

- Detect chemical reactions

- Forces of particle $i$

$$\vec{f_i} = m_i \vec{a_i} = -\nabla_i U(t) \qquad (1)$$

- Potential

$$U(t) = \frac{1}{2} \sum_{i=1}^{N} \sum_{\substack{j=1 \\ j \neq i}}^{N} U_{i,j}(\|\vec{r}_{i,j}\|) \qquad (2)$$

**Algorithm 1** Overview of the main Molecular Dynamics routine.

1: **for** $i = 1$ to $M$ **do**
2:    $t \leftarrow t + dt$
3:    compute_forces($\vec{r}, \vec{f}$)
4:    integrate($\vec{r}, \vec{f}, \vec{v}, dt$)
5:    //Do something with the data
6: **end for**

- *compute_forces* has a complexity of $\mathcal{O}(N^2)$

**Algorithm 2** Compute_forces routine.

1:  **for** $i = 1$ to $N$ **do**
2:      $\vec{f_i} \leftarrow 0$
3:      **for** $j = 1$ to $N$ **do**
4:          $\vec{r}_{i,j} \leftarrow \vec{r_j} - \vec{r_i}$
5:          $f_{i,j} \leftarrow \text{compute\_force}(\|\vec{r}_{i,j}\|)$
6:          $\vec{f_i} \leftarrow \vec{f_i} + f_{i,j}\ \vec{r}_{i,j}$
7:      **end for**
8:  **end for**

---

**Algorithm 3** Naive OpenACC compute_forces routine.

1: #pragma acc kernels
2: **for** $i = 1$ to $N$ **do**
3: $\quad \vec{f_i} \leftarrow 0$
4: $\quad$ **for** $j = 1$ to $N$ **do**
5: $\quad\quad \vec{r_{i,j}} \leftarrow \vec{r_j} - \vec{r_i}$
6: $\quad\quad f_{i,j} \leftarrow \text{compute\_force}(\|\vec{r_{i,j}}\|)$
7: $\quad\quad \vec{f_i} \leftarrow \vec{f_i} + f_{i,j} \ \vec{r_{i,j}}$
8: $\quad$ **end for**
9: **end for**

---

- Inner loop can not be parallelized
  - Loop-carried dependencies

**Algorithm 4** Improved compute_forces routine.

```
 1: #pragma acc kernels
 2: for i = 1 to N do
 3:     f_i ← 0
 4:     #pragma acc loop reduction(+:f_i)
 5:     for j = 1 to N do
 6:         r_{i,j} ← r_j − r_i
 7:         f_{i,j} ← compute_force(||r_{i,j}||)
 8:         f_i ← f_i + f_{i,j} r_{i,j}
 9:     end for
10: end for
```
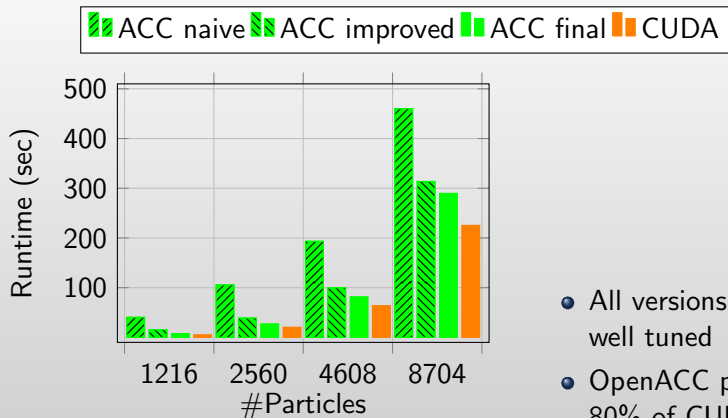
- Good: Inner loop can be parallelized
- Bad: Arrays are reallocated in every iteration

**Algorithm 5** Overview of the main Molecular Dynamics routine with an OpenACC data region.

1: #pragma acc data create ($\vec{r}$[0:N],$\vec{f}$[0:N])
2: **for** $i = 1$ to $M$ **do**
3:    $t \leftarrow t + dt$
4:    compute_forces($\vec{r}, \vec{f}$)
5:    integrate($\vec{r}, \vec{f}, \vec{v}, dt$)
6:    //Do something with the data
7: **end for**

**Algorithm 6** Final compute_forces routine.

1: #pragma acc update device($\vec{r}$[0:N])
2: #pragma acc kernels present($\vec{r}$[0:N], $\vec{f}$[0:N])
3: **for** $i = 1$ to $N$ **do**
4:    $\vec{f_i} \leftarrow 0$
5:    #pragma acc loop reduction($+$:$\vec{f_i}$)
6:    **for** $j = 1$ to $N$ **do**
7:       $\vec{r_{i,j}} \leftarrow \vec{r_j} - \vec{r_i}$
8:       $f_{i,j} \leftarrow$ compute_force($\|\vec{r_{i,j}}\|$)
9:       $\vec{f_i} \leftarrow \vec{f_i} + f_{i,j}\ \vec{r_{i,j}}$
10:    **end for**
11: **end for**
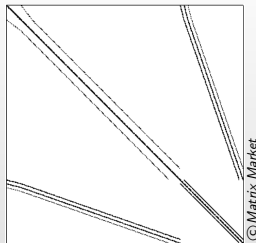12: #pragma acc update host($\vec{f}$[0:N])

Figure: Runtime of a Molecular Dynamics (MD) Simulation for different problem sizes over 10.000 iterations. All calculations are run in double precision. OpenMP: 16 core SMP node. OpenACC/Cuda: Nvidia Quadro 6000 GPU.
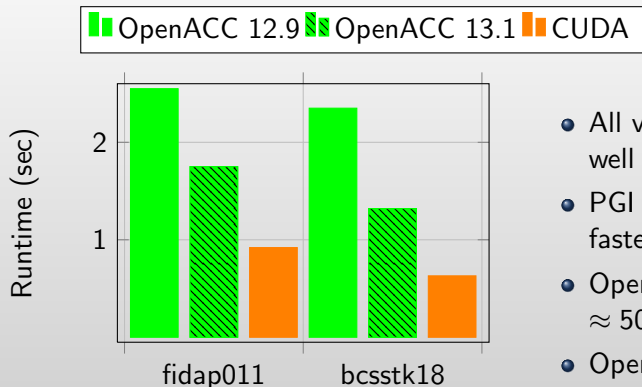
- All versions are equally well tuned
- OpenACC performs at 80% of CUDA

Conjugate Gradient Method

# Conjugate Gradient Method

© Matrix Market

- Iterative solver

- Solve a large sparse linear system

$$Ax = b \qquad (3)$$

- Frequently arise from *partial differential equations* in physics

- Runtime dominated by *Sparse Matrix-Vector Multiplication* SPMV

# Performance

Figure: Runtime of a Conjugate Gradient (CG) Method for two sparse matrices. All calculations are run in double precision. OpenMP: 16 core SMP node. OpenACC/Cuda: Nvidia Quadro 6000 GPU.

- All versions are equally well tuned
- PGI 13.1 50%/80% faster than PGI 12.9
- OpenACC performs at $\approx 50\%$ of CUDA
- OpenMP outperforms CUDA

# Productivity

**RWTH**AACHEN
UNIVERSITY

- Function calls require inlining

- PGI compiler does not support C++

- Limited debugging support for PGI compiler
  - Revert to debugging the logic of your application

|    | OpenMP | OpenACC | CUDA |
|----|--------|---------|------|
| MD | 23     | 16      | 92   |
| CG | 8      | 16      | 156  |

Table: Number of added and modified lines of source code for each case study and paradigm with respect to the serial version.

- Few added/modified lines of source code
- Data transfers are straight forward
- Reductions require almost no additional effort
- No need to worry about "boundary conditions"
- Compiler is able to tune for a specific coprocessor

- High productivity (if you don't run into compiler bugs)

- Decent performance

- Limited debugging support for PGI compiler

- Makes coprocessor programming more straight forward
  - C code $\rightarrow$ OpenACC code $\rightarrow$ CUDA code

- High productivity (if you don't run into compiler bugs)

- Decent performance

- Limited debugging support for PGI compiler

- Makes coprocessor programming more straight forward
  - C code $\rightarrow$ OpenACC code $\rightarrow$ CUDA code

## Thank you for your attention.