# High-Performance Computing for Stencil Computations Using a High-Level Domain-Specific Language
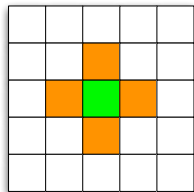
**Justin Holewinski**      Tom Henretty      Kevin Stock
Louis-Noël Pouchet      Atanas Rountev      P. Sadayappan

Department of Computer Science and Engineering
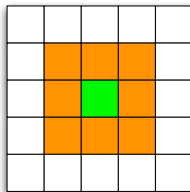The Ohio State University

Presented at WOLFHPC 2011

May 31, 2011

# Stencil Computations


5-pt 2D Stencil


9-pt 2D Stencil

```
1  for(i = 1; i < N-1; ++i) {
2    for(j = 1; j < N-1; ++j) {
3      A[i][j] = CNST * (B[i  ][j  ] +
4                        B[i  ][j-1] +
5                        B[i  ][j+1] +
6                        B[i-1][j  ] +
7                        B[i+1][j  ]);
8    }
9  }
```

▶ Operate on each point in a discrete *n*-dimensional space

▶ Use neighboring points in computation

▶ Often surrounded by time loop

▶ Have diverse boundary conditions

# Domain-Specific Language for Stencils

Why do we need a domain-specific language?

# Domain-Specific Language for Stencils

Why do we need a domain-specific language?

- ▶ Easier for application developers and scientists
  - ▶ Write stencil as *point-function* and *grid* instead of loop nest

# Domain-Specific Language for Stencils

Why do we need a domain-specific language?

- Easier for application developers and scientists
  - Write stencil as *point-function* and *grid* instead of loop nest

- More opportunity for compiler optimization
  - Restricted to a simple expression language
  - Not restricted by C/C++/Fortran specification
    e.g. aliasing, memory life-cycle
  - Control-flow is implicit instead of discovered at compile-time
  - Iteration domain is easily obtained, enabling polyhedral
    transformations for tiling, parallelism, memory optimizations
  - Computations on grids ease dependency analysis
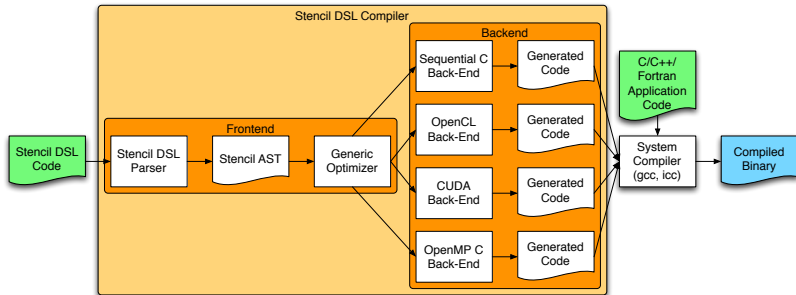
# Domain-Specific Language for Stencils

Why do we need a domain-specific language?

- ▶ Easier for application developers and scientists
    - ▶ Write stencil as *point-function* and *grid* instead of loop nest

- ▶ More opportunity for compiler optimization
    - ▶ Restricted to a simple expression language
    - ▶ Not restricted by C/C++/Fortran specification
      e.g. aliasing, memory life-cycle
    - ▶ Control-flow is implicit instead of discovered at compile-time
    - ▶ Iteration domain is easily obtained, enabling polyhedral
      transformations for tiling, parallelism, memory optimizations
    - ▶ Computations on grids ease dependency analysis

### Goal

*Use high-level abstractions to achieve write-once performance
portability for stencil computations.*

# Domain-Specific Language for Stencils



Stencil Compiler Workflow

# Domain-Specific Language for Stencils

Define stencil operation as a *point-function* over a *grid* using
`[time]grid[i-offset][j-offset]` notation:

```
1 pointfunction five_point_avg(p) {
2   float ONE_FIFTH;
3   ONE_FIFTH  = 0.2;
4   [1]p[0][0] = ONE_FIFTH*([0]p[-1][0] + [0]p[0][-1] + [0]p[0][0]
5                                       + [0]p[0][1]  + [0]p[1][0]);
6 }
```

# Domain-Specific Language for Stencils

Define stencil operation as a *point-function* over a *grid* using
`[time]grid[i-offset][j-offset]` notation:

```
1  pointfunction five_point_avg(p) {
2    float ONE_FIFTH;
3    ONE_FIFTH  = 0.2;
4    [1]p[0][0] = ONE_FIFTH*([0]p[-1][0] + [0]p[0][-1] + [0]p[0][0]
5                                        + [0]p[0][1]  + [0]p[1][0]);
6  }
```

Define stencil range, functions, and convergence:

```
1  iterate 1000 {
2    stencil jacobi_2d {
3      [0][0:Nx-1]    : [1]a[0][0] = [0]a[0][0];
4      [Ny-1][0:Nx-1] : [1]a[0][0] = [0]a[0][0];
5      [0:Ny-1][0]    : [1]a[0][0] = [0]a[0][0];
6      [0:Ny-1][Nx-1] : [1]a[0][0] = [0]a[0][0];
7  
8      [1:Ny-2][1:Nx-2] : five_point_avg(a);
9    }
10 
11   reduction max_diff max {
12     [0:Ny-1][0:Nx-1] : [1]a[0][0] - [0]a[0][0];
13   }
14 } check (max_diff < .00001) every 4 iterations
```

# Domain-Specific Language for Stencils

```
1  int  Nx;
2  int  Ny;
3  grid g [Ny][Nx];
4
5  float griddata a on g at 0,1;
6
7  pointfunction five_point_avg(p) {
8    float ONE_FIFTH;
9    ONE_FIFTH  = 0.2;
10   [1]p[0][0] = ONE_FIFTH*([0]p[-1][0] + [0]p[0][-1] + [0]p[0][0]
11                                       + [0]p[0][1]  + [0]p[1][0]);
12 }
13
14 iterate 1000 {
15   stencil jacobi_2d {
16     [0][0:Nx-1]    : [1]a[0][0] = [0]a[0][0];
17     [Ny-1][0:Nx-1] : [1]a[0][0] = [0]a[0][0];
18     [0:Ny-1][0]    : [1]a[0][0] = [0]a[0][0];
19     [0:Ny-1][Nx-1] : [1]a[0][0] = [0]a[0][0];
20
21     [1:Ny-2][1:Nx-2] : five_point_avg(a);
22   }
23
24   reduction max_diff max {
25     [0:Ny-1][0:Nx-1] : [1]a[0][0] - [0]a[0][0];
26   }
27 } check (max_diff < .00001) every 4 iterations
```

Complete stencil program

Floating-Point Throughput
- ▸ Need fine-grain and coarse-grain parallelism

# Optimizing CPU vs. GPU Performance

Floating-Point Throughput
- ▶ Need fine-grain and coarse-grain parallelism
- ▶ On a CPU
  - ▶ Use vector processing units (SIMD)
  - ▶ Use threads to exploit multi-/many-cores

# Optimizing CPU vs. GPU Performance

Floating-Point Throughput

- ▶ Need fine-grain and coarse-grain parallelism
- ▶ On a CPU
    - ▶ Use vector processing units (SIMD)
    - ▶ Use threads to exploit multi-/many-cores
- ▶ On a GPU
    - ▶ Exploit SIMT parallelism across hundreds of cores
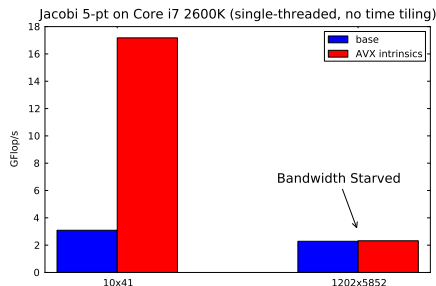    - ▶ Multiprocessors operate in lock-step ⇒ divergence = BAD

*But this is not the whole story...*

# Optimizing CPU vs. GPU Performance

Floating-Point Throughput

- Need fine-grain and coarse-grain parallelism
- On a CPU
  - Use vector processing units (SIMD)
  - Use threads to exploit multi-/many-cores
- On a GPU
  - Exploit SIMT parallelism across hundreds of cores
  - Multiprocessors operate in lock-step $\Rightarrow$ divergence = BAD

*But this is not the whole story...*



Jacobi 5-pt on Core i7 2600K (single-threaded, no time tiling)

# Optimizing CPU vs. GPU Performance

Memory Hierarchy

- ▶ Increasingly complex (multiple levels)
- ▶ Fast but small on-chip memory
- ▶ Slow but abundant off-chip memory

# Optimizing CPU vs. GPU Performance

Memory Hierarchy

- Increasingly complex (multiple levels)
- Fast but small on-chip memory
- Slow but abundant off-chip memory
- On a CPU
    - Exploit hardware caches through data re-use

# Optimizing CPU vs. GPU Performance

Memory Hierarchy

- ► Increasingly complex (multiple levels)
- ► Fast but small on-chip memory
- ► Slow but abundant off-chip memory
- ► On a CPU
  - ► Exploit hardware caches through data re-use
- ► On a GPU
  - ► Exploit per-multiprocessor shared/local memory
  - ► Maximize work per read/write operation
- ► *Need time tiling to efficiently utilize available main memory bandwidth*

# Optimizing Stencils on GPUs

Typical Approach

- ► Use spatial tiling to distribute work among thread blocks
- ► Use shared/local memory as program-controlled cache

# Optimizing Stencils on GPUs

Typical Approach

- ▶ Use spatial tiling to distribute work among thread blocks
- ▶ Use shared/local memory as program-controlled cache

Problems

- ▶ Global (off-chip) memory latency is high
- ▶ Limited data re-use within a thread block
- ▶ Cannot schedule enough threads to hide memory latency
- ▶ Traditional time tiling is not efficient due to branch divergence and a lack of memory access coalescing

# Optimizing Stencils on GPUs

Typical Approach

- Use spatial tiling to distribute work among thread blocks
- Use shared/local memory as program-controlled cache

Problems

- Global (off-chip) memory latency is high
- Limited data re-use within a thread block
- Cannot schedule enough threads to hide memory latency
- Traditional time tiling is not efficient due to branch divergence and a lack of memory access coalescing

Result

- Compute units are mostly idle waiting for memory operations to complete

# Optimizing Stencils on GPUs

Typical Approach
- ► Use spatial tiling to distribute work among thread blocks
- ► Use shared/local memory as program-controlled cache

Problems
- ► Global (off-chip) memory latency is high
- ► Limited data re-use within a thread block
- ► Cannot schedule enough threads to hide memory latency
- ► Traditional time tiling is not efficient due to branch divergence and a lack of memory access coalescing

Result
- ► Compute units are mostly idle waiting for memory operations to complete

A possible solution?
- ► *Overlapped tiling*

# Overlapped Tiling

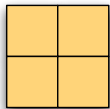Replace inter-tile communication with redundant computation

- ► Tile borders are redundantly computed by all neighboring tiles
- ► Trades extra FLOPs for a decrease in needed synchronization
- ► Enables time tiling without skewing (introduces divergence, load imbalance, and bank conflicts)

# Overlapped Tiling

Replace inter-tile communication with redundant computation

- ► Tile borders are redundantly computed by all neighboring tiles
- ► Trades extra FLOPs for a decrease in needed synchronization
- ► Enables time tiling without skewing (introduces divergence, load imbalance, and bank conflicts)

Originally proposed by Krishnamoorthy et al. for parallelization

- ► We want fully-automatic code generation for arbitrary stencils
- ► Use OpenCL for performance-portable code generation, but tune parameters for different GPU architectures

# Overlapped Tiling

Replace inter-tile communication with redundant computation

- ▶ Tile borders are redundantly computed by all neighboring tiles
- ▶ Trades extra FLOPs for a decrease in needed synchronization
- ▶ Enables time tiling without skewing (introduces divergence, load imbalance, and bank conflicts)

Originally proposed by Krishnamoorthy et al. for parallelization

- ▶ We want fully-automatic code generation for arbitrary stencils
- ▶ Use OpenCL for performance-portable code generation, but tune parameters for different GPU architectures

Let us look at an example for a $2 \times 2$ tile with a time tile size of 2...

Tile at time t +1

# Overlapped Tiling



Computed in time step t

Data needed at time t +1

# Overlapped Tiling
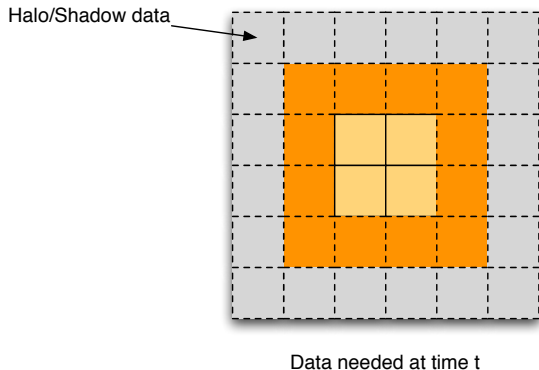
Also computed by neighboring tiles



Computation at time t

Halo/Shadow data

Data needed at time t

# Overlapped Tiling

GPU Implementation

- ▶ Schedule extra threads for redundant border cell computations
  - ▶ In general, need $(n + 2 * r * (t - 1)) \times (m + 2 * r * (t - 1))$ threads

# Overlapped Tiling

GPU Implementation

- ▶ Schedule extra threads for redundant border cell computations
  - ▶ In general, need $(n + 2 * r * (t - 1)) \times (m + 2 * r * (t - 1))$ threads
- ▶ Use shared memory to store results across time
  - ▶ Only need to access global memory in first and last time step of tile

# Overlapped Tiling

GPU Implementation

- ▶ Schedule extra threads for redundant border cell computations
  - ▶ In general, need $(n + 2 * r * (t - 1)) \times (m + 2 * r * (t - 1))$ threads

- ▶ Use shared memory to store results across time
  - ▶ Only need to access global memory in first and last time step of tile

- ▶ Synchronize threads, not blocks, after each time step
  - ▶ Thread synchronization efficiently supported in hardware; block synchronization is not

# Overlapped Tiling

GPU Implementation

- ▶ Schedule extra threads for redundant border cell computations
  - ▶ In general, need $(n + 2 * r * (t - 1)) \times (m + 2 * r * (t - 1))$ threads

- ▶ Use shared memory to store results across time
  - ▶ Only need to access global memory in first and last time step of tile

- ▶ Synchronize threads, not blocks, after each time step
  - ▶ Thread synchronization efficiently supported in hardware; block synchronization is not

- ▶ Use host to synchronize across time tiles

```
1  for(t = 0; t < TIME_STEPS; t += TIME_TILE_SIZE) {
2    invoke_kernel(input, output);
3    swap(input, output);
4    // Implicit barrier
5  }
```

# What about block size?

Block size considerations

- ▶ Block size has large impact on performance
- ▶ Need enough threads to keep compute units busy...
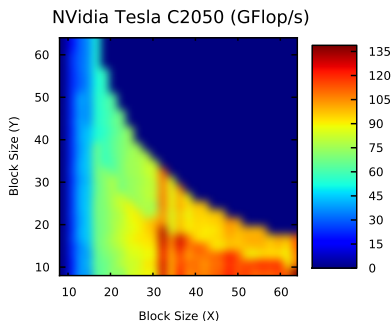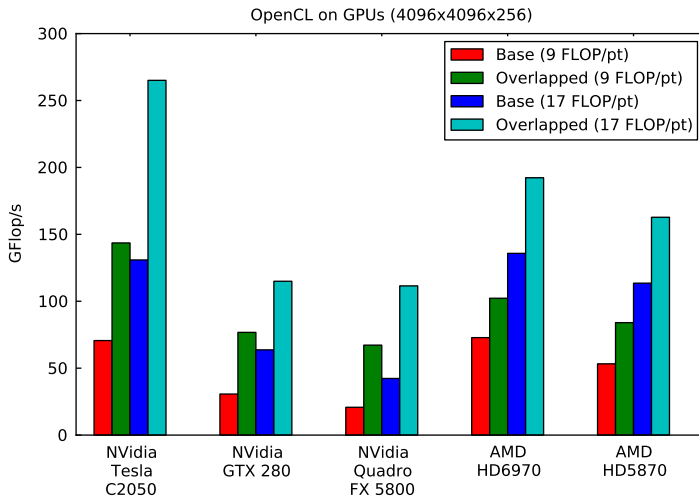
# What about block size?

Block size considerations

- ▶ Block size has large impact on performance
- ▶ Need enough threads to keep compute units busy...
- ▶ ... but it is also beneficial to use smaller blocks to increase the number of available registers per block

# What about block size?

Block size considerations

- ▶ Block size has large impact on performance
- ▶ Need enough threads to keep compute units busy...
- ▶ ... but it is also beneficial to use smaller blocks to increase the number of available registers per block
- ▶ Problem size: $4096 \times 4096 \times 256$



NVidia Tesla C2050 (GFlop/s)

AMD Radeon HD 6970 (GFlop/s)

# Arithmetic Intensity

Arithmetic intensity matters too...



OpenCL on GPUs (4096x4096x256)

# Performance



Jacobi 5-pt on Multi-Core with OpenMP

Jacobi 5-pt on GPU with OpenCL

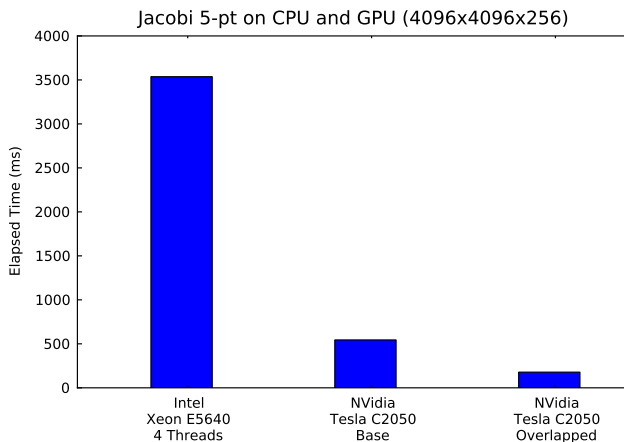- ▶ Fixed CPU tile sizes
- ▶ Fixed GPU block/tile sizes

# Conclusion

- A DSL for stencils enables high productivity and performance
  - Higher-level for application developers
  - More information for compilers
  - Increased performance-portability

- Overlapped tiling enables high-performance stencils on GPUs
  - Trade redundant computation for less communication
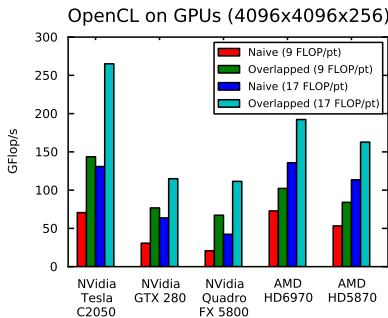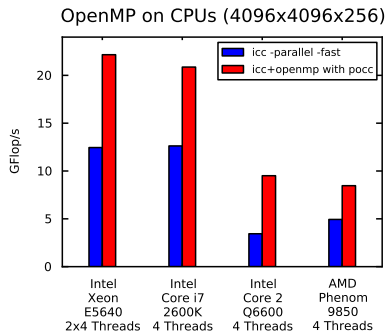  - Exploit high compute-per-memory-op ratio on GPUs
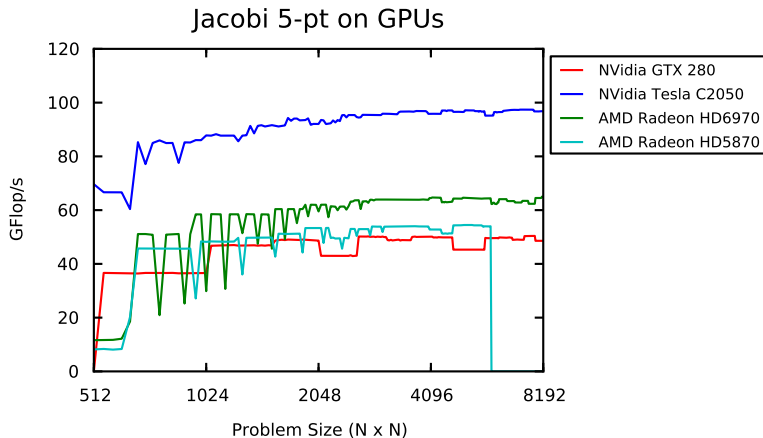
# Questions?

# Performance Evaluation



Jacobi 5-pt on CPU and GPU (4096x4096x256)

GPU Block Size: $64 \times 8$ (512 of 1024 max)

# Performance Evaluation



FP Through-put for Jacobi 9-pt

# Performance Evaluation



Jacobi 5-pt on GPUs

Problem Size Evaluation for GPUs