# Optimizing Stencil Computations

## March 18, 2013

THE UNIVERSITY OF UTAH

---

# Administrative

THE UNIVERSITY OF UTAH

---

# Stencil Computations

A stencil defines the value of a grid point in a d-dimensional spatial grid at time t as a function of neighboring grid points at recent times before t.

**(a)** 1D 5-Point Stencil      **(b)** 2D 5-Point Stencil

THE UNIVERSITY OF UTAH

---

# Stencil Computations, Performance Issues

• Bytes per flop ratio is O(1)
• Most machines cannot supply data at this rate, leading to memory bound computation
• Some reuse, but difficult to exploit fully, and interacts with parallelization

How to maximize performance:
• Avoid extraneous memory traffic, such as cache capacity/conflict misses
• Bandwidth optimizations to maximize utility of memory transfers
• Maximize in-core performance

07/20/11

THE UNIVERSITY OF UTAH

## Learn More: StencilProbe

- See http://people.csail.mit.edu/skamil/projects/stencilprobe/
- Several variations of Heat Equation, to be discussed
- Can instantiate to measure performance impact

07/20/11

## Example: Heat Equation

```
for (t=0; t<timesteps; t++)  {  // time step loop
 for (k=1; k<nz-1; k++) {
   for (j=1; j<ny-1; j++) {
     for (i=1; i<nx-1; i++) {
       // 3-d 7-point stencil
       B[i][j][k] = A[i][j][k+1] + A[i][j][k-1] +
         A[i][j+1][k] + A[i][j-1][k] + A[i+1][j][k] +
         A[i-1][j][k] – 6.0 * A[i][j][k] / (fac*fac);
     }
   }
 }
 temp_ptr = A;
 A = B;
 B = temp_ptr;
}
```

What if nx, ny, nz large?

07/20/11

## Heat Equation, Add Tiling

```
for (t=0; t<timesteps; t++)  {  // time step loop
 for (jj = 1; jj < ny-1; jj+=TJ) {
   for (ii = 1; ii < nx - 1; ii+=TI) {
     for (k=1; k<nz-1; k++) {
       for (j = jj; j < MIN(jj+TJ,ny - 1); j++) {
         for (i = ii; i < MIN(ii+TI,nx - 1); i++) {
           // 3-d 7-point stencil
           B[i][j][k] = A[i][j][k+1] + A[i][j][k-1] +
             A[i][j+1][k] + A[i][j-1][k] + A[i+1][j][k] +
             A[i-1][j][k] – 6.0 * A[i][j][k] / (fac*fac);
         }
       }
     }
   }
 temp_ptr = A;
 A = B;
 B = temp_ptr;
}
```
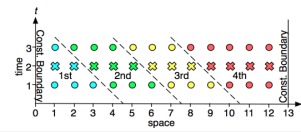
Note the reuse across time steps!

07/20/11

## Heat Equation, Time Skewing

```
for (kk=1; kk < nz-1; kk+=tz) {
 for (jj = 1; jj < ny-1; jj+=ty) {
   for (ii = 1; ii < nx - 1; ii+=tx) {
     for (t=0; t<timesteps; t++)  {  // time step loop
       … calculate bounds from t and slope …
       for (k=blockMin_z; k < blockMax_z; k++) {
         for (j=blockMin_y; j < blockMax_y; j++) {
           for (i=blockMin_x; i < blockMax_x; i++) {
             // 3-d 7-point stencil
             B[i][j][k] = A[i][j][k+1] + A[i][j][k-1] +
               A[i][j+1][k] + A[i][j-1][k] + A[i+1][j][k] +
               A[i-1][j][k] – 6.0 * A[i][j][k] / (fac*fac);
           }
         }
       }
     }
 temp_ptr = A;
 A = B;
 B = temp_ptr;
}
```
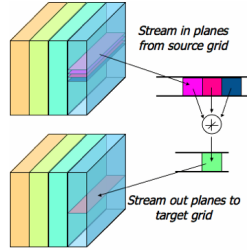


07/20/11

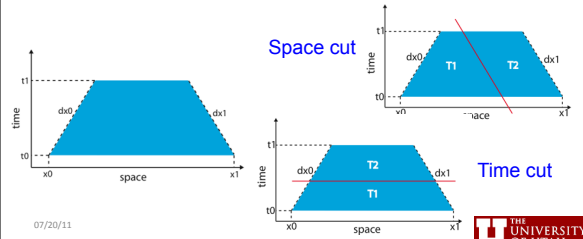# Heat Equation, Circular Queue

- See probe_heat_circqueue.c



Stream in planes from source grid

Stream out planes to target grid

# Heat Equation, Cache Oblivious

- See probe_heat_oblivious.c
- Idea: Recursive decomposition to cutoff point
- Implicit tiling: in both space and time
- Simpler code than complex tiling, but introduces overhead
- Encapsulated in Pochoir DSL (next slide)



Space cut

Time cut

# Example Pochoir Stencil Compiler Specification

```
1   #define mod(r,m)  ((r)%(m) + ((r)<0)? (m):0)

2   Pochoir_Boundary_2D(heat_bv, a, t, x, y)
3     return a.get(t,mod(x,a.size(1)),mod(y,a.size(0)));
4   Pochoir_Boundary_End

5   int main(void) {

6     Pochoir_Shape_2D 2D_five_pt[] = {{1,0,0}, {0,1,0},
        {0,-1,0}, {0,-1,-1}, {0,0,-1}, {0,0,1}};
7     Pochoir_2D heat(2D_five_pt);

8     Pochoir_Array_2D(double) u(X, Y);
9     u.Register_Boundary(heat_bv);
10    heat.Register_Array(u);

11    Pochoir_Kernel_2D(heat_fn, t, x, y)
12      u(t+1, x, y) = CX * (u(t, x+1, y) - 2 * u(t, x,
          y) + u(t, x-1, y) + CY * (u(t, x, y+1) - 2
          * u(t, x, y) + u(t, x, y-1)) + u(t, x, y);
13    Pochoir_Kernel_End

14    for (int x = 0; x < X; ++x)
15      for (int y = 0; y < Y; ++y)
16        u(0, x, y) = rand();

17    heat.Run(T, heat_fn);

18    for (int x = 0; x < X; ++x)
19      for (int y = 0; y < Y; ++y)
20        cout << u(T, x, y);

22    return 0;
23  }
```
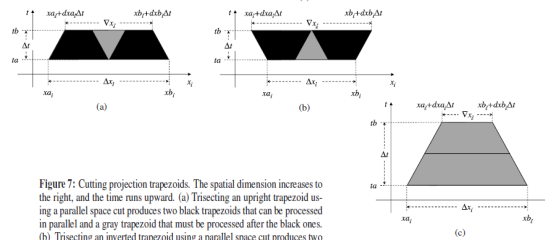
# Parallel Stencils in Pochoir



Figure 7: Cutting projection trapezoids. The spatial dimension increases to the right, and the time runs upward. (a) Trisecting an upright trapezoid using a parallel space cut produces two black trapezoids that can be processed in parallel and a gray trapezoid that must be processed after the black ones. (b) Trisecting an inverted trapezoid using a parallel space cut produces two black trapezoids that can be processed in parallel and a gray trapezoid that must be processed before the black ones. (c) A time cut produces a lower and an upper trapezoid where the lower trapezoid must be processed before the upper.
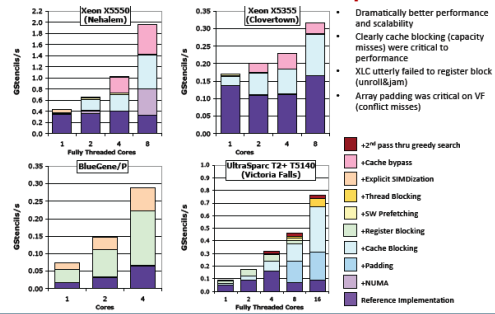
## General Approach to Parallel Stencils

• Always safe to parallelize within a time step
• Circular queue and time skewing encapsulate "tiles" that are independent

---

## Results for Heat Equation



• Dramatically better performance and scalability
• Clearly cache blocking (capacity misses) were critical to performance
• XLC utterly failed to register block (unroll&jam)
• Array padding was critical on VF (conflict misses)

Legend:
+2nd pass thru greedy search
+Cache bypass
+Explicit SIMDization
+Thread Blocking
+SW Prefetching
+Register Blocking
+Cache Blocking
+Padding
+NUMA
Reference Implementation

Reference: K. Datta et al., *"Stencil Computation Optimization and Autotuning on State-of-the-Art Multi-core Architectures"*, SC08.
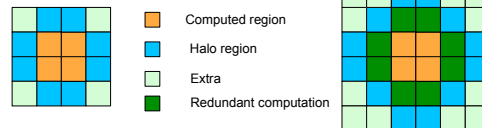
07/20/11

---

## What about GPUs?

• Two recent papers:
"Auto-Generation and Auto-Tuning of 3D Stencil Codes on GPU Clusters," Zhang and Mueller, CGO 2012.

"High-Performance Code Generation for Stencil Computations on GPU Architectures," Holewinski et al., ICS 2012.

• Key issues:
—Exploit reuse in shared memory.
—Avoid fetching from global memory.
—Thread decomposition to support global memory coalescing.

---

## Overlapped Tiling for Halo Regions (or Ghost Zones)

• Input data exceeds output result (as in Sobel)
• Halo region or ghost zone extends the per-thread data decomposition to encompass additional input
• An (n+2)x(n+2) halo region is needed to compute an nxn block if subscript expressions are of the form ± 1, for example
• By expanding the halo region, we can trade off redundant computation for reduced accesses to global memory when parallelizing across time steps.



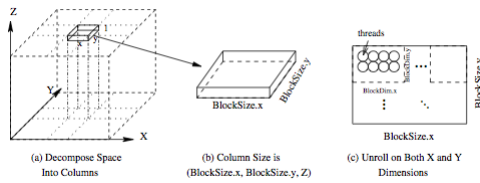Computed region
Halo region
Extra
Redundant computation

2-d 5-point stencil example

## 2.5D Decomposition

• Partition such that each thread block sweeps over the z-axis and processes one plane at a time.



(a) Decompose Space Into Columns

(b) Column Size is (BlockSize.x, BlockSize.y, Z)

(c) Unroll on Both X and Y Dimensions

---

## Resulting Code



(a) With Corner Accesses

(b) Without Corner Accesses

**Figure 5: Stencil Kernel Templates**

---

## Other Optimizations

• X dimension delivers coalesced global memory accesses
  • Pad to multiples of 32 stencil elements
• Halo regions are aligned to 128-bit boundaries
• Input (parameter) arrays are padded to match halo region, to share indexing.
• BlockSize.x is maximized to avoid non-coalesced accesses to halo region
• Blocks are square to reduce area of redundancy.
• Use of shared memory for input.
• Use of texture fetch for input.

---

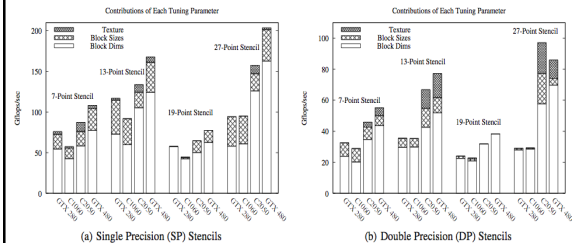## Performance Across Stencils and Architectures



(a) Single Precision (SP) Stencils

(b) Double Precision (DP) Stencils

**Figure 8: Stencil Tuning Effect Breakups**
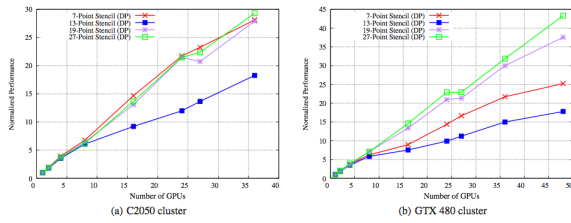
# GPU Cluster Performance



(a) C2050 cluster

(b) GTX 480 cluster

Figure 11: Weak Scaling of DP Stencils on GPU Clusters