

L5: Writing Correct Programs, cont.

CS6963

1
L5: Writing Correct Programs

Outline

- How to tell if your parallelization is correct?
- Race conditions and data dependences
- Tools that detect race conditions
- Abstractions for writing correct parallel code
 - Mapping these abstractions to CUDA
- Reading (if you can find it):
 - "Optimizing Compilers for Modern Architectures: A Dependence-Based Approach", Allen and Kennedy, 2002, Ch. 2.

CS6963

2
L5: Writing Correct Programs

Administrative

- Next assignment (a homework) given out on Monday

CS6963

3
L5: Writing Correct Programs

Is this CUDA code correct?

```

__host__ callkernel() {
    dim3 blocks{10};
    dim3 threads{100};
    float *d_array;
    ...
    cudaMalloc(&d_array1,...);
    cudaMalloc(&d_array2,...);
    kernelcode<<<blocks,threads,
    0>>>{d_array1, d_array2,
    1000};
}

__global__ kernelcode(float * d_array1,
    d_array2, int N) {
    float result;
    for (int i=0; i<N; i++) {
        d_array1[threadIdx] +=
        d_array2[blockIdx][i];
    }
    for (int i=1; i<N; i++) {
        result += d_array1[threadIdx-1];
    }
}

```

CS6963

4
L5: Writing Correct Programs

Threads Access Shared Memory!

- Global memory and shared memory within an SM can be freely accessed by multiple threads
- Requires appropriate sequencing of memory accesses across threads to same location **if at least one access is a write**
 - Recall using `__syncthreads()` within a thread block for synchronization
 - Not to be used for different blocks within a grid

CS6963

5
LS: Writing Correct Programs

Is this CUDA code correct?

```

__host__ callkernel() {
    dim3 blocks{10};
    dim3 threads{100};
    float *d_array;
    ...
    cudaMalloc(&d_array1,...);
    cudaMalloc(&d_array2,...);
    kernelcode<<<blocks,threads,
    0>>>{d_array1, d_array2,
    1000};
}

__global kernelcode(float * d_array1,
    d_array2, int N) {
    for (int i=0; i<N; i++) {
        d_array1[threadIdx] +=
        d_array2[blockIdx][i];
    }
    __syncthreads();
    for (int i=1; i<N; i++) {
        result += d_array1[threadIdx-1];
    }
}

```

CS6963

6
LS: Writing Correct Programs

More Formally: Race Condition or Data Dependence

- A **race condition** exists when the result of an execution depends on the **timing** of two or more events.
- A **data dependence** is an ordering on a pair of memory operations that must be preserved to maintain correctness.

CS6963

7
LS: Writing Correct Programs

How about other Shared Memory Architectures?

- Race detection software (e.g., Intel ThreadChecker)
 - Trace memory accesses for each thread
 - Compare addresses accessed by each thread
 - Race condition exists if, between synchronization points,
 - multiple threads access the same memory location
 - and, at least one access is a write

CS6963

8
LS: Writing Correct Programs

What can we do to debug parallelization in CUDA?

- deviceemu code (to be emulated on host)
 - Support for pthread debugging?
- Can compare GPU output to CPU output, or compare GPU output to device emulation output
 - Race condition may still be present
- Or can (try to) prevent introduction of race conditions (remainder of lecture)

CS6963

9

LS: Writing Correct Programs



Data Dependence

- **Definition:**
Two memory accesses are involved in a data dependence if they may refer to the same memory location and one of the references is a write.

A data dependence can either be between two distinct program statements or two different dynamic executions of the same program statement.
- Two important uses of data dependence information (among others):
Parallelization: no data dependence between two computations → parallel execution safe
Locality optimization: absence of data dependences & presence of reuse → reorder memory accesses for better data locality (next week)

CS6963

10

LS: Writing Correct Programs



Data Dependence of Scalar Variables

True (flow) dependence

$$a = a$$

Anti-dependence

$$a = a$$

Output dependence

$$a = a$$

Input dependence (for locality)

$$= a$$

Definition: Data dependence exists from a reference instance i to i' iff either i or i' is a write operation and i and i' refer to the same variable i executes before i'

CS6963

11

LS: Writing Correct Programs



Some Definitions (from Allen & Kennedy)

- **Definition 2.5:**
 - Two computations are equivalent if, on the same inputs,
 - they produce identical outputs
 - the outputs are executed in the same order
- **Definition 2.6:**
 - A reordering transformation
 - changes the order of statement execution
 - without adding or deleting any statement executions.
- **Definition 2.7:**
 - A reordering transformation preserves a dependence if
 - it preserves the relative execution order of the dependences' source and sink.

CS6963

12

LS: Writing Correct Programs



Fundamental Theorem of Dependence

- **Theorem 2.2:**
 - Any reordering transformation that preserves every dependence in a program preserves the meaning of that program.
- Now we will discuss abstractions and algorithms to determine whether reordering transformations preserve dependences...

CS6963

13
LS: Writing Correct Programs

Parallelization as a Reordering Transformation in CUDA

```

__host__ callkernel() {
    dim3 blocks{bx,by};
    dim3 threads{tx,ty,tz};
    ...
    kernelcode<<<blocks,threadsO>>><<
    args>>;
}
__global__ kernelcode{<args>} {
    /* code refers to threadIdx.x,
    threadIdx.y, threadIdx.z, blockDim.x,
    blockDim.y */
}

__host__ callkernel() {
    for (int bldx_x=0; bldx_x<bx; bldx_x++) {
    for (int bldx_y=0; bldx_y<by; bldx_y++) {
    for (int tldx_x=0; tldx_x<tx; tldx_x++) {
    for (int tldx_y=0; tldx_y<ty; tldx_y++) {
    for (int tldx_z=0; tldx_z<tz; tldx_z++) {
    for (int tldx_z=0; tldx_z<tz; tldx_z++) {
    /* code refers to tldx_x, tldx_y, tldx_z,
    bldx_x, bldx_y */
    }
    }
    }
    }
    }
}

```

EQUIVALENT?

CS6963

14
LS: Writing Correct Programs

In Today's Lecture: Parallelizable Loops

Forall (or CUDA kernels or Doall) loops:
Loops whose iterations can execute in parallel (a particular reordering transformation)

Example

```
forall (i=1; i<=n; i++)
    A[i] = B[i] + C[i];
```

Meaning?

Each iteration can execute independently of others
Free to schedule iterations in any order

Why are parallelizable loops an important concept for data-parallel programming models?

CS6963

15
LS: Writing Correct Programs

CUDA Equivalent to "Forall"

```

__host__ callkernel() {
    forall (int bldx_x=0; bldx_x<bx; bldx_x++) {
    forall (int bldx_y=0; bldx_y<by; bldx_y++) {
    forall (int tldx_x=0; tldx_x<tx; tldx_x++) {
    forall (int tldx_y=0; tldx_y<ty; tldx_y++) {
    forall (int tldx_z=0; tldx_z<tz; tldx_z++) {

    /* code refers to tldx_x, tldx_y, tldx_z,
    bldx_x, bldx_y */
    }
    }
    }
    }
    }
}

```

CS6963

16
LS: Writing Correct Programs

Data Dependence for Arrays

```
for (i=2; i<5; i++)
  A[i] = A[i-2]+1;

for (i=1; i<=3; i++)
  A[i] = A[i]+1;
```

Loop-Carried dependence

Loop-Independent dependence

- Recognizing parallel loops (intuitively)
 - Find data dependences in loop
 - No dependences crossing iteration boundary → parallelization of loop's iterations is safe

CS6963
17
LS: Writing Correct Programs

1. Characterize Iteration Space

```
for (i=1; i<=5; i++)
  for (j=i; j<=7; j++)
    ...
```

1 ≤ i ≤ 5
i ≤ j ≤ 7

- **Iteration instance:** represented as coordinates in iteration space
 - n -dimensional discrete cartesian space for n deep loop nests
- **Lexicographic order:** Sequential execution order of iterations
 $[1,1], [1,2], \dots, [1,6], [1,7], [2,2], [2,3], \dots, [2,6], \dots$
- Iteration I (a vector) is lexicographically less than I' , $I < I'$, iff there exists $c (i_1, \dots, i_{c-1}) = (i'_1, \dots, i'_{c-1})$ and $i_c < i'_c$.

CS6963
18
LS: Writing Correct Programs

2. Compare Memory Accesses across Dynamic Instances in Iteration Space

```
N = 6;
for (i=1; i<N; i++)
  for (j=1; j<N; j++)
    A[i+1,j+1] = A[i,j] * 2.0;
```

I=[1,1], Write A[2,2]

I'=[2,2], Read A[2,2]

How to describe relationship between two dynamic instances?
e.g., $I=[1,1]$ and $I'=[2,2]$

CS6963
19
LS: Writing Correct Programs

Distance Vectors

```
N = 6;
for (i=1; i<N; i++)
  for (j=1; j<N; j++)
    A[i+1,j+1] = A[i,j] * 2.0;
```

Distance vector = $[1,1]$

- A loop has a distance vector D if there exists data dependence from iteration vector I to a later vector I' , and $D = I' - I$.
- Since $I' > I$, $D \geq 0$. (D is lexicographically greater than or equal to 0).

CS6963
20
LS: Writing Correct Programs

Distance and Direction Vectors

- Distance vectors: (infinitely large set)

$$\left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ -1 \end{bmatrix}, \dots \right)$$

- Direction vectors: (realizable if 0 or lexicographically positive)

$$([=,=], [=,<], [<,>], [<,<], [<,<])$$

- Common notation:

0 =
+ <
- >
+/- *

CS6963

21
LS: Writing Correct Programs

Parallelization Test: 1-Dimensional Loop

- Examples:

```
for (j=1; j<N; j++)
  A[j] = A[j] + 1;
```

```
for (j=1; j<N; j++)
  B[j] = B[j-1] + 1;
```

- Dependence (Distance and Direction) Vectors?

- Test for parallelization:

- A loop is parallelizable if for all data dependences $D \in \mathbf{D}$,
 $D = 0$

CS6963

22
LS: Writing Correct Programs

n-Dimensional Loop Nests

```
for (i=1; i<=N; i++)
  for (j=1; j<=N; j++)
    A[i][j] = A[i][j-1]+1;
```

```
for (i=1; i<=N; i++)
  for (j=1; j<=N; j++)
    A[i][j] = A[i-1][j+1]+1;
```

- Distance and direction vectors?

- Definition:**

$D = (d_1, \dots, d_n)$ is loop-carried at level i if d_i is the first nonzero element.

CS6963

23
LS: Writing Correct Programs

A Few Words about n- Dimensional Arrays in C

- Largely conceptual, due to difficulty in expressing this in C for dynamically allocated data structures
- Imagine the following macros,

```
#define 2dAccess(i,j,dim_i) \
  i+j*dim_i
```

```
#define 3dAccess(i,j,k,dim_i,dim_j) \
  i+j*dim_i + k*dim_i*dim_j
```

CS6963

24
LS: Writing Correct Programs

Test for Parallelization

The i th loop of an n -dimensional loop is parallelizable if there does not exist any level i data dependences.

The i th loop is parallelizable if for all dependences

$$D = (d_1, \dots, d_n),$$

either

$$(d_1, \dots, d_{i-1}) > 0$$

or

$$(d_1, \dots, d_i) = 0$$

CS6963

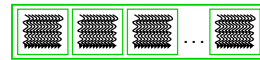
25
L5: Writing Correct Programs



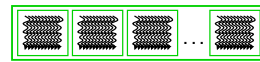
Safe Parallelization of CUDA Code

- Dependences must be carried by
 - (a) Loops that execute on the host
 - OR, loops that execute within a kernel function

*May be able to use synchronization for dependences across threads, but not across blocks (subtle distinction)



(a) Dependences carried by host code



```
__global compute() {
  for (j=-1; j<n; j++)
    A[threadIdx[j]] =
      A[threadIdx[j]-1];
}
```

(b) Dependence carried within thread code

CS6963

L5: Writing Correct Programs



Parallelization Algorithm

- For each pair of dynamic accesses to the same array within a loop nest:
 - determine if there exists a dependence between that pair
- Key points:
 - n^2 tests for n accesses in loop!
 - a single access is compared with itself
 - includes accesses in all loops within a nest

CS6963

27
L5: Writing Correct Programs



Dependence Testing

- Question so far:
 - What is the distance/direction (in the iteration space) between two dynamic accesses to the same memory location?
- Simpler question:
 - Can two data accesses ever refer to the same memory location?

```
for (i=11; i<=20; i++)      for (i=11; i<=20; i++)
  A[i] = A[i-1] + 3;        A[i] = A[i-10] + 1;
```

CS6963

28
L5: Writing Correct Programs



Restrict to an Affine Domain

```
for (i=1; i<N; i++)
  for (j=1; j<N; j++) {
    A[i+2*j+3, 4*i+2*j, 3*i] = ...;
    ... = A[1, 2*i+1, j];
  }
```

- Only use loop bounds and array indices which are integer linear functions of loop variables.

- Non-affine example:**

```
for (i=1; i<N; i++)
  for (j=1; j<N; j++) {
    A[i*j] = A[i*(j-1)];
    A[B[i]] = A[B[j]];
  }
```

29
LS: Writing Correct Programs



Equivalence to Integer Programming

- Need to determine if $F(i) = G(i')$, where i and i' are iteration vectors, with constraints $i, i' \geq L, U \geq i, i'$

- Example:**

```
for (i=1; i<=100; i++)
  A[i] = A[i-1];
```

- Inequalities:**

$1 \leq iw \leq 100, \quad ir = iw - 1, \quad ir \leq 100$
integer vector $I, \quad AI \leq b$

- Integer Programming is NP-complete**

- $O(\text{size of the coefficients})$
- $O(n^r)$

CS6963

30
LS: Writing Correct Programs



Example

$0 \leq iw \leq 100, ir = iw - 1, ir \leq 100$

Solve

$$\begin{bmatrix} -10 \\ 10 \\ -11 \\ 1-1 \\ 01 \end{bmatrix} \begin{bmatrix} iw \\ ir \end{bmatrix} \leq \begin{bmatrix} 1 \\ 100 \\ -1 \\ 1 \\ 100 \end{bmatrix} \quad \text{Solution exist?}$$

Replace instances of iw with ir
Eliminate ir
Simplifies to $0 < 99$ (yes, solution exists!)

CS6963

31
LS: Writing Correct Programs



Introducing Omega Calculator

- A software tool used in compilers
- Description:
 - Solves "Presburger formulas", logical formulas built from affine constraints over integer variables, logical connectives and quantifiers.
 - Can formulate these dependence problems, and derive existence of dependences and distances.
- Relevant to locality optimizations as well, next week's focus
- Can download from:
 - <http://www.cs.utah.edu/~chunchen/omega/>
- Also available from CADE Linux machines in:
 - `~cs6963/bin/oc`

CS6963

32
LS: Writing Correct Programs



Using Omega Calculator

- **Example:**

```
for (i=2; i<=100; i++)
  A[i] = A[i-1];
```
 - **Define relation** $iw = i$, and $iw = ir-1$ in the iteration space $2 \leq i \leq 100$.

```
R := {[iw] -> [ir] :
  2 <= iw, ir <= 100      /* iteration space */
  && iw < ir              /* looking for loop-carried dep */
  && iw = ir-1;           /* can they be the same? */
```
- R := {[iw] -> [ir] : 2 <= iw, ir <= 100 && iw < ir && iw = ir - 1};

Result: {[iw] -> [iw+1] : 2 <= iw <= 99}

CS6963

33
LS: Writing Correct Programs

Using Omega Calculator, cont.

- **Example:**

```
for (i=20; i<=29; i++)
  A[i] = A[i-10];
```
 - **Define relation** $iw = i$, and $iw = ir-10$ in the iteration space $20 \leq i \leq 29$.

```
R := {[iw] -> [ir] :
  20 <= iw, ir <= 29      /* iteration space */
  && iw < ir              /* looking for loop-carried dep */
  && iw = ir-10;         /* can they be the same? */
```
- R := {[iw] -> [ir] : 20 <= iw, ir <= 29 && iw < ir && iw = ir - 10};

Result: {[iw] -> [ir] : FALSE }

CS6963

34
LS: Writing Correct Programs

2-D Omega Example

- **Example:**

```
for (i=0; i<n; i++)
  for (j=i; j<n; j++)
    a[i][j+1] = a[n][j];
```
 - **Formula (more complex):**

```
R := {[iw,jw] -> [ir,jr] : exists(n : /* unbound variable */
  1 <= iw <= n && iw <= jw <= n /* iteration space */
  && 1 <= ir <= jr <= n /* loop-carried dependence? */
  && jw+1 = jr && ir = n); /* access expression */
```
- Result:** {[iw,jw] -> [ir,jr] : FALSE }

CS6963

35
LS: Writing Correct Programs

Calculating Distance Vectors

- **Example from before:**

```
for (i=1; i<=100; i++)
  for (j=1; j<=100; j++)
    A[i][j] = A[i-1][j+1]+1;
```
 - **Omega formula:**

```
R := {[iw, jw] -> [di, dj] : exists (ir, jr : /* read iters unbound */
  1 <= iw, ir <= 100 && 1 <= jw, jr <= 100 /* iteration space */
  && iw = ir - 1 && jw = jr + 1 && /* access exprs */
  && di = ir - iw && dj = jr - jw); /* distances */
```
- Result:** {[iw,jw] -> [1,-1] : 1 <= iw <= 99 && 2 <= jw <= 100}

CS6963

36
LS: Writing Correct Programs

Aside: What about dependences for other data structures?

- Pointer-based
 - Pointer alias analysis
 - Shape analysis
- Objects
 - Escape analysis
- In practice
 - Lots of `#pragma` and special flags for programmer to assert no dependence

CS6963

37
LS: Writing Correct Programs

Homework Assigned Monday

- Example questions
 - Given some sequential code, do the following:
 - Show distance vectors and decide which loop can be parallelized
 - Show Omega dependence relations
 - Show correct CUDA code
 - Memory hierarchy optimization
 - Simple tiling example
 - Identify safety of code transformations
 - Given description of goal, show CUDA code to manage memory hierarchy

CS6963

38
LS: Writing Correct Programs

Summary of Lecture

- Data dependence can be used to determine the safety of reordering transformations such as parallelization
 - preserving dependences = preserving "meaning"
- Iteration, distance and direction vectors are abstractions for understanding whether reordering transformations preserve dependences.
 - Parallelization of CUDA kernel programs can be viewed as a reordering transformation of a sequential implementation
- Dependence testing on array accesses in loops has been shown to be equivalent to integer programming.
 - Omega calculator demonstrated

CS6963

39
LS: Writing Correct Programs

What's Ahead

- Next week
 - Homework assignment on Monday
 - Managing the memory hierarchy
 - Initial discussion of projects
- February 16: President's Day holiday
- February 18:
 - Jim Guilkey to present MPM for projects
- February 20 (Friday):
 - Make up class (we'll discuss control flow)

CS6963

40
LS: Writing Correct Programs