

CS4961 Parallel Programming

Lecture 7: Introduction to SIMD

Mary Hall
September 14, 2010

09/14/2010

CS4961

Homework 2, Due Friday, Sept. 10, 11:59 PM

- To submit your homework:
 - Submit a PDF file
 - Use the "handin" program on the CADE machines
 - Use the following command:


```
"handin cs4961 hw2 <prob2file>"
```

Problem 1 (based on #1 in text on p. 59):

Consider the Try2 algorithm for "count3s" from Figure 1.9 of p.19 of the text. Assume you have an input array of 1024 elements, 4 threads, and that the input data is evenly split among the four processors so that accesses to the input array are local and have unit cost. Assume there is an even distribution of appearances of 3 in the elements assigned to each thread which is a constant we call NTPT. What is a bound for the memory cost for a particular thread predicted by the CTA expressed in terms of λ and NTPT.

09/14/2010

CS4961



Homework 2, cont

Problem 2 (based on #2 in text on p. 59), cont.:

Now provide a bound for the memory cost for a particular thread predicted by CTA for the Try4 algorithm of Fig. 114 on p. 23 (or Try3 assuming each element is placed on a separate cache line).

Problem 3:

For these examples, how is algorithm selection impacted by the value of NTPT?

Problem 4 (in general, not specific to this problem):

How is algorithm selection impacted by the value of λ ?

09/14/2010

CS4961



Programming Assignment 1 Due Wednesday, Sept. 21 at 11:59PM

- Logistics:
 - You'll use water.eng.utah.edu (a Sun UltraSparc T2), for which all of you have accounts that match the userid and password of your CADE Linux account.
 - Compile using "cc -O3 -xopenmp p01.c"
- Write the prefix sum computation from HW1 in OpenMP using the test harness found on the website.
 - What is the parallel speedup of your code as reported by the test harness?
 - If your code does not speed up, you will need to adjust the parallelism granularity, the amount of work each processor does between synchronization points. You can adjust this by changing numbers of threads, and frequency of synchronization. You may also want to think about reducing the parallelism overhead, as the solutions we have discussed introduce a lot of overhead.
 - What happens when you try different numbers of threads or different schedules?

09/14/2010

CS4961



Programming Assignment 1, cont.

- What to turn in:
 - Your source code so we can see your solution
 - A README file that describes at least three variations on the implementation or parameters and the performance impact of those variations.
 - handin "cs4961 p1 <gzipped tarfile>"
- Lab hours:
 - Thursday afternoon and Tuesday afternoon

09/14/2010

CS4961



Review: Predominant Parallel Control Mechanisms

Name	Meaning	Examples
Single Instruction, Multiple Data (SIMD)	A single thread of control, same computation applied across "vector" elts	Array notation as in Fortran 90: <code>A[1:n] = A[1:n] + B[1:n]</code>
Multiple Instruction, Multiple Data (MIMD)	Multiple threads of control, processors periodically synch	Parallel loop: <code>forall (i=0; i<n; i++)</code>
Single Program, Multiple Data (SPMD)	Multiple threads of control, but each processor executes same code	Processor-specific code: <code>if (\$myid == 0) { } }</code>

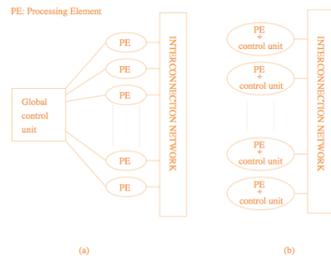
09/01/2009

CS4961

6



SIMD and MIMD Architectures: What's the Difference?



A typical SIMD architecture (a) and a typical MIMD architecture (b).

Slide source: Grama et al., Introduction to Parallel Computing, <http://www.users.cs.umn.edu/~karypis/parbook>

09/01/2009

CS4961

7



Overview of SIMD Programming

- Vector architectures
- Early examples of SIMD supercomputers
- TODAY Mostly
 - Multimedia extensions such as SSE and AltiVec
 - Graphics and games processors
 - Accelerators (e.g., ClearSpeed)
- Is there a dominant SIMD programming model
 - Unfortunately, NO!!!
- Why not?
 - Vector architectures were programmed by scientists
 - Multimedia extension architectures are programmed by systems programmers (almost assembly language!)
 - GPUs are programmed by games developers (domain-specific libraries)

09/08/2009

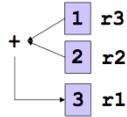
CS4961

8

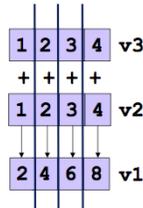


Scalar vs. SIMD in Multimedia Extensions

Scalar: `add r1,r2,r3`



SIMD: `vadd<sws> v1,v2,v3`



09/08/2009

CS4961

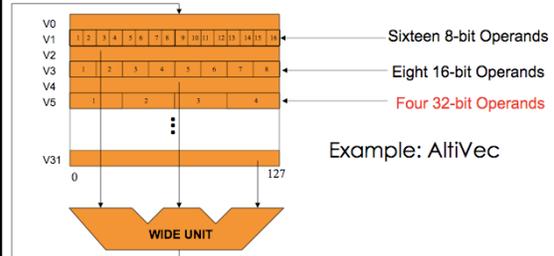
9



Multimedia Extension Architectures

• At the core of multimedia extensions

- SIMD parallelism
- Variable-sized data fields:
- Vector length = register width / type size



Example: Altivec

09/08/2009

CS4961

10



Multimedia / Scientific Applications

- Image
 - Graphics : 3D games, movies
 - Image recognition
 - Video encoding/decoding : JPEG, MPEG4
- Sound
 - Encoding/decoding: IP phone, MP3
 - Speech recognition
 - Digital signal processing: Cell phones
- Scientific applications
 - Double precision Matrix-Matrix multiplication (DGEMM)
 - $Y[] = a * X[] + Y[]$ (SAXPY)

09/10/2010

CS4961

11



Characteristics of Multimedia Applications

- Regular data access pattern
 - Data items are contiguous in memory
- Short data types
 - 8, 16, 32 bits
- Data streaming through a series of processing stages
 - Some temporal reuse for such data streams
- Sometimes ...
 - Many constants
 - Short iteration counts
 - Requires saturation arithmetic

09/10/2010

CS4961

12



Why SIMD

- +More parallelism
 - +When parallelism is abundant
 - +SIMD in addition to ILP
- +Simple design
 - +Replicated functional units
- +Small die area
 - +No heavily ported register files
 - +Die area: +MAX-2(HP): 0.1% +VIS(Sun): 3.0%
- Must be explicitly exposed to the hardware
 - By the compiler or by the programmer

09/08/2009

CS4961

13



Programming Multimedia Extensions

- Language extension
 - Programming interface similar to function call
 - C: built-in functions, Fortran: intrinsics
 - Most native compilers support their own multimedia extensions
 - GCC: `-faltivec, -msse2`
 - Altivec: `dst= vec_add(src1, src2);`
 - SSE2: `dst= _mm_add_ps(src1, src2);`
 - BG/L: `dst= __fpadd(src1, src2);`
 - No Standard !
- Need automatic compilation

09/08/2009

CS4961

14



Programming Complexity Issues

- High level: Use compiler
 - may not always be successful
- Low level: Use intrinsics or inline assembly tedious and error prone
- Data must be aligned, and adjacent in memory
 - Unaligned data may produce incorrect results
 - May need to copy to get adjacency (overhead)
- Control flow introduces complexity and inefficiency
- Exceptions may be masked

09/08/2009

CS4961

15



1. Independent ALU Ops

```
R = R + XR * 1.08327
G = G + XG * 1.89234
B = B + XB * 1.29835
```



R	=	R	+	XR	*	1.08327
G	=	G	+	XG	*	1.89234
B	=	B	+	XB	*	1.29835

09/10/2010

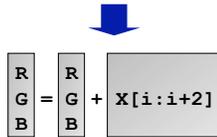
CS4961

16



2. Adjacent Memory References

```
R = R + X[i+0]
G = G + X[i+1]
B = B + X[i+2]
```



09/10/2010

CS4961

17



3. Vectorizable Loops

```
for (i=0; i<100; i+=1)
  A[i+0] = A[i+0] + B[i+0]
```

09/10/2010

CS4961

18



3. Vectorizable Loops

```
for (i=0; i<100; i+=4)
  A[i+0] = A[i+0] + B[i+0]
  A[i+1] = A[i+1] + B[i+1]
  A[i+2] = A[i+2] + B[i+2]
  A[i+3] = A[i+3] + B[i+3]
```

```
for (i=0; i<100; i+=4)
  A[i:i+3] = B[i:i+3] + C[i:i+3]
```

09/10/2010

CS4961

19



4. Partially Vectorizable Loops

```
for (i=0; i<16; i+=1)
  L = A[i+0] - B[i+0]
  D = D + abs(L)
```

09/10/2010

CS4961

20



4. Partially Vectorizable Loops

```
for (i=0; i<16; i+=2)
  L = A[i+0] - B[i+0]
  D = D + abs(L)
  L = A[i+1] - B[i+1]
  D = D + abs(L)
```



```
for (i=0; i<16; i+=2)
  L0 = A[i:i+1] - B[i:i+1]
  L1 = A[i:i+1] - B[i:i+1]
  D = D + abs(L0)
  D = D + abs(L1)
```

09/10/2010

CS4961

21



Exploiting SLP with SIMD Execution

- Benefit:
 - Multiple ALU ops → One SIMD op
 - Multiple ld/st ops → One wide mem op
- Cost:
 - Packing and unpacking
 - Reshuffling within a register
 - Alignment overhead

09/10/2010

CS4961

22



Packing/Unpacking Costs

```
C = A + 2
D = B + 3
```

```
C = A + 2
D = B + 3
```

09/10/2010

CS4961

23



Packing/Unpacking Costs

- Packing source operands
 - Copying into contiguous memory

```
A = f()
B = g()
C = A + 2
D = B + 3
```

```
A = f()
B = g()
C = A + 2
D = B + 3
```

09/10/2010

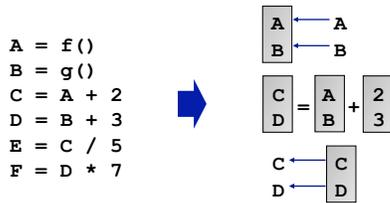
CS4961

24



Packing/Unpacking Costs

- Packing source operands
 - Copying into contiguous memory
- Unpacking destination operands
 - Copying back to location



09/10/2010

CS4961

25

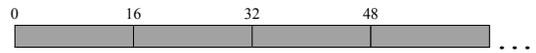


Alignment Code Generation

- Aligned memory access
 - The address is always a multiple of 16 bytes
 - Just one superword load or store instruction

```

float a[64];
for (i=0; i<64; i+=4)
    Va = a[i:i+3];
    
```



09/10/2010

CS4961

26



Alignment Code Generation (cont.)

- Misaligned memory access
 - The address is always a non-zero constant offset away from the 16 byte boundaries.
 - Static alignment: For a misaligned load, issue two adjacent aligned loads followed by a merge.

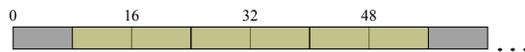
```

float a[64];
for (i=0; i<60; i+=4)
    Va = a[i+2:i+5];
    
```

➔

```

float a[64];
for (i=0; i<60; i+=4)
    V1 = a[i:i+3];
    V2 = a[i+4:i+7];
    Va = merge(V1, V2, 8);
    
```



09/10/2010

CS4961

27



- Statically align loop iterations

```

float a[64];
for (i=0; i<60; i+=4)
    Va = a[i+2:i+5];

float a[64];
Sa2 = a[2]; Sa3 = a[3];
for (i=2; i<62; i+=4)
    Va = a[i+2:i+5];
    
```

09/10/2010

CS4961

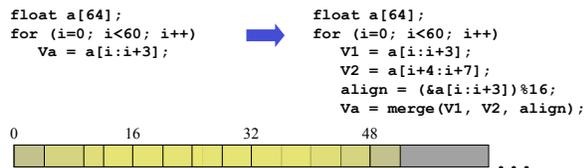
28



Alignment Code Generation (cont.)

- **Unaligned** memory access

- The offset from 16 byte boundaries is varying or not enough information is available.
- Dynamic alignment: The merging point is computed during run time.



09/10/2010

CS4961

29



SIMD in the Presence of Control Flow

```
for (i=0; i<16; i++)
  if (a[i] != 0)
    b[i]++;
```



```
for (i=0; i<16; i+=4){
  pred = a[i:i+3] != (0, 0, 0, 0);
  old = b[i:i+3];
  new = old + (1, 1, 1, 1);
  b[i:i+3] = SELECT(old, new, pred);
}
```

Overhead:

Both control flow paths are always executed !

09/10/2010

CS4961

30



An Optimization: Branch-On-Superword-Condition-Code



```
for (i=0; i<16; i+=4){
  pred = a[i:i+3] != (0, 0, 0, 0);
  branch-on-none(pred) L1;
  old = b[i:i+3];
  new = old + (1, 1, 1, 1);
  b[i:i+3] = SELECT(old, new, pred);
L1:
}
```

09/10/2010

CS4961

31



Control Flow

- Not likely to be supported in today's commercial compilers
 - Increases complexity of compiler
 - Potential for slowdown
 - Performance is dependent on input data
- Many are of the opinion that SIMD is not a good programming model when there is control flow.
- But speedups are possible!

09/10/2010

CS4961

32



Nuts and Bolts

- What does a piece of code really look like?

```
for (i=0; i<100; i+=4)
  A[i:i+3] = B[i:i+3] + C[i:i+3]
```

```
for (i=0; i<100; i+=4) {
  __m128 btmp = _mm_load_ps(float B[i]);
  __m128 ctmp = _mm_load_ps(float C[i]);
  __m128 atmp = _mm_add_ps(__m128 btmp, __m128 ctmp);
  void_mm_store_ps(float A[i], __m128 atmp);
}
```

09/10/2010

CS4961

33



Wouldn't you rather use a compiler?

- Intel compiler is pretty good
 - `icc -mssse3 -vecreport3 <file.c>`
- Get feedback on why loops were not "vectorized"
- First programming assignment
 - Use compiler and rewrite code examples to improve vectorization
 - One example: write in low-level intrinsics

09/10/2010

CS4961

34



Next Time

- Discuss Red-Blue computation, problem 10 on page 111 (not assigned, just to discuss)
- More on Data Parallel Algorithms

09/14/2010

CS4961

