CMPT 880/479 Multicore Programming and Architecture CMPT 880/479 Multicore Programming and Architecture

Unit 1: Performance Tuning

CMPT 880/479 Multicore Programming and Architecture

Unit 1: Performance Tuning

Today's Agenda

- Why all this talk about sequential performance?
 - Seldom taught
 - Even more rarely practiced in coursework
 - Key foundation for efficient parallel code

Today's Agenda

- Why all this talk about sequential performance?
 - Seldom taught
 - Even more rarely practiced in coursework
 - Key foundation for efficient parallel code

Today's Agenda

- Why all this talk about sequential performance?
 - Seldom taught
 - Even more rarely practiced in coursework
 - Key foundation for efficient parallel code

- Lecture on sequential performance tuning
 - Overview, much of it should be review
 - Help you tune code
 - Help you understand odd performance effects

The Process of Performance Tuning

- Reduce resource usage to "go faster"
 - **Runtime**: wall clock time for a single task
 - Throughput: items of work per unit time

- Reduce resource usage to "go faster"
 - **Runtime**: wall clock time for a single task
 - Throughput: items of work per unit time
- Memory usage

- Reduce resource usage to "go faster"
 - **Runtime**: wall clock time for a single task
 - Throughput: items of work per unit time
- Memory usage
- Or, energy consumption
 - Same computation, less energy
 - Why? battery life, power & cooling costs, fan noise
 - Run fast and sleep --or-- run slowly just to meet deadlines

- Reduce resource usage to "go faster"
 - **Runtime**: wall clock time for a single task
 - Throughput: items of work per unit time
- Memory usage

• Or, energy consumption

- Same computation, less energy
- Why? battery life, power & cooling costs, fan noise
- Run fast and sleep --or-- run slowly just to meet deadlines
- Either way: **software tuning basically the same**
 - In fact, processors dynamically trade energy and performance
 - Dynamic voltage/frequency scaling, Core i7 "turbo mode"

• Starting point: always some pile of existing code

- Starting point: always some pile of existing code
- **Refactoring:** No change to visible external behavior

- Starting point: always some pile of existing code
- **Refactoring:** No change to visible external behavior
- Add new functionality
 - Change behavior

- Starting point: always some pile of existing code
- **Refactoring:** No change to visible external behavior
- Add new functionality
 - Change behavior

• **Performance optimization**: empirical & experimental

- 1. wait
- 2. benchmark and profile, identify bottleneck
- 3. Modify code, test functionality, measure perf., revert if not beneficial
- Ongoing performance regression testing

- Starting point: always some pile of existing code
- **Refactoring:** No change to visible external behavior
- Add new functionality
 - Change behavior
- **Performance optimization**: empirical & experimental
 - 1. wait
 - 2. benchmark and profile, identify bottleneck
 - 3. Modify code, test functionality, measure perf., revert if not beneficial
 - Ongoing performance regression testing
- "We should forget about small efficiencies, say about 97% of the time: premature optimization is the root of all evil"
 Donald Knuth

• Where is the code spending its time?

- Need representative input (full sized, long running)
- Various profiling tools (gprof, Valgrind's cachegrind, Intel's VTune)
- Look beyond code profile (collect by object, algorithm change, etc.)

• Where is the code spending its time?

- Need representative input (full sized, long running)
- Various profiling tools (gprof, Valgrind's cachegrind, Intel's VTune)
- Look beyond code profile (collect by object, algorithm change, etc.)

• What is the cause of the slow performance?

• Might not always be what you think

• Where is the code spending its time?

- Need representative input (full sized, long running)
- Various profiling tools (gprof, Valgrind's cachegrind, Intel's VTune)
- Look beyond code profile (collect by object, algorithm change, etc.)

• What is the cause of the slow performance?

• Might not always be what you think

• Hardware performance counters

- Various profiling tools collect statistics from the hardware
- Examples: stall cycles, cache miss rates, branch mispredictions, etc.
- Can even associate them with specific instructions
- Using sampling, can be done with low overhead (a few percent)

• Where is the code spending its time?

- Need representative input (full sized, long running)
- Various profiling tools (gprof, Valgrind's cachegrind, Intel's VTune)
- Look beyond code profile (collect by object, algorithm change, etc.)

• What is the cause of the slow performance?

• Might not always be what you think

• Hardware performance counters

- Various profiling tools collect statistics from the hardware
- Examples: stall cycles, cache miss rates, branch mispredictions, etc.
- Can even associate them with specific instructions
- Using sampling, can be done with low overhead (a few percent)
- Focused experiments (example: large v. small data set)

Performance Tuning of Sequential Code

• Fast parallel code is fast sequential code running in parallel

- Fast parallel code is fast sequential code running in parallel
- Runtime = # of inst. * cycles per inst. * clock period
 - Want to minimize all three

- Fast parallel code is fast sequential code running in parallel
- Runtime = # of inst. * cycles per inst. * clock period
 - Want to minimize all three
- When tuning code, consider:
 - "Dynamic" instruction count
 - Micro-architectural impacts
 - Instruction level parallelism of code, branch mispredictions
 - Memory hierarchy
 - Cache misses, etc.

- Dynamic instruction count is first-order performance metric
 - "Dynamic" means at runtime

- Dynamic instruction count is first-order performance metric
 - "Dynamic" means at runtime

• First order model: runtime is ~ to instruction count

- Assumes all instruction take a single cycle
- Con: ignores all micro-architectural and memory effects
- Pro: architecture independent, easier to reason about, captures most important metric

- Dynamic instruction count is first-order performance metric
 - "Dynamic" means at runtime

• First order model: runtime is ~ to instruction count

- Assumes all instruction take a single cycle
- Con: ignores all micro-architectural and memory effects
- Pro: architecture independent, easier to reason about, captures most important metric

• Reducing instruction count

- Better algorithm or approach to problem
- Avoid software bloat in terms of too much layering, too general
- Write code the compiler can easily optimize (and turn on -O3)
 - Compilers are really smart, but lack higher level context

Cost of Function Calls

Cost of Function Calls

• Function calls add overhead

• Obvious overhead: "Call" and "Return" instruction overheads

Cost of Function Calls

• Function calls add overhead

• Obvious overhead: "Call" and "Return" instruction overheads

• Less obvious

- Putting parameters/arguments into right registers
 - Or on in-memory stack
- Calling into a dynamically linked library can be a bit more expensive
- Virtual functions (C++ and Java): adds a table lookup on call path
Cost of Function Calls

• Function calls add overhead

• Obvious overhead: "Call" and "Return" instruction overheads

• Less obvious

- Putting parameters/arguments into right registers
 - Or on in-memory stack
- Calling into a dynamically linked library can be a bit more expensive
- Virtual functions (C++ and Java): adds a table lookup on call path

• Worse yet, breaks compiler optimization boundaries

- Compilers today typically don't look at code in other C files
- After a call, compiler doesn't know what has changed or not
- Compilers can inline functions, especially in C header files

Cost of Function Calls

• Function calls add overhead

• Obvious overhead: "Call" and "Return" instruction overheads

• Less obvious

- Putting parameters/arguments into right registers
 - Or on in-memory stack
- Calling into a dynamically linked library can be a bit more expensive
- Virtual functions (C++ and Java): adds a table lookup on call path

• Worse yet, breaks compiler optimization boundaries

- Compilers today typically don't look at code in other C files
- After a call, compiler doesn't know what has changed or not
- Compilers can inline functions, especially in C header files
- Aside: "extern inline" function in .h file is a fast as a macro

- Instruction latencies
 - One cycle: integer adds, subtract, compare, shift, etc.
 - Few cycles: loads/stores (best case), int multiply, floating point
 - Dozens of cycles: integer divide (89 cycles for 64-bit on Core i7)

- Instruction latencies
 - One cycle: integer adds, subtract, compare, shift, etc.
 - Few cycles: loads/stores (best case), int multiply, floating point
 - Dozens of cycles: integer divide (89 cycles for 64-bit on Core i7)
- Non-instructions (on most machines)
 - Sqrt, exponentiation, sine/cosine, modulo (%), etc.
 - Trend toward more hardware support for these operations

- Instruction latencies
 - One cycle: integer adds, subtract, compare, shift, etc.
 - Few cycles: loads/stores (best case), int multiply, floating point
 - Dozens of cycles: integer divide (89 cycles for 64-bit on Core i7)
- Non-instructions (on most machines)
 - Sqrt, exponentiation, sine/cosine, modulo (%), etc.
 - Trend toward more hardware support for these operations
- Compilers can by help
 - Selects cheapest instructions (strength reduction)
 - Replace multiply by constant power-of-two with shift: "x * 2" with "x << 1"
 - Tries to eliminate redundant computation

- Superscalar execution: three or four insn per cycle max
 - Restrictions on instruction mix (different function unit types)
 - Loads/stores are most limited (just one or maybe two per cycle)
 - Compiler reorganizes and schedules code with this in mind

- Superscalar execution: three or four insn per cycle max
 - Restrictions on instruction mix (different function unit types)
 - Loads/stores are most limited (just one or maybe two per cycle)
 - Compiler reorganizes and schedules code with this in mind
- Out-of-order (dynamic scheduling)
 - Code runs faster, but makes performance difficult to understand

- Superscalar execution: three or four insn per cycle max
 - Restrictions on instruction mix (different function unit types)
 - Loads/stores are most limited (just one or maybe two per cycle)
 - Compiler reorganizes and schedules code with this in mind
- Out-of-order (dynamic scheduling)
 - Code runs faster, but makes performance difficult to understand
- Branching
 - "Taken" branches disrupt fetch
 - Branch mis-predictions flush pipeline, throws away work
 - Cost hard to quantify, but often dozens of cycles
 - Try to avoid unpredictable branches

- Superscalar execution: three or four insn per cycle max
 - Restrictions on instruction mix (different function unit types)
 - Loads/stores are most limited (just one or maybe two per cycle)
 - Compiler reorganizes and schedules code with this in mind
- Out-of-order (dynamic scheduling)
 - Code runs faster, but makes performance difficult to understand
- Branching
 - "Taken" branches disrupt fetch
 - Branch mis-predictions flush pipeline, throws away work
 - Cost hard to quantify, but often dozens of cycles
 - Try to avoid unpredictable branches
- Loop accelerators (if loop body is <N insn, goes faster)

```
int distance(int a, int b)
{
   return (a > b) ? a-b : b-a;
}
```

```
int distance(int a, int b)
{
    return (a > b) ? a-b : b-a;
}
    gcc -m32 -O3
```

distance: pushl %ebp movl %esp, %ebp movl 8(%ebp), %edx movl 12(%ebp), %eax cmpl %eax, %edx jg .L6 subl %edx, %eax popl %ebp ret .L6: subl %eax, %edx movl %edx, %eax popl %ebp

ret

```
int distance(int a, int b)
{
   return (a > b) ? a-b : b-a;
}
```

<u>gcc -m32 -O3</u> <u>gcc -m32 -march=core2 -O3</u>

distance:

distance:

pushl	% ebp	pushl	%ebp
movl	%esp, %ebp	movl	%esp, %ebp
movl	8(%ebp), %edx	pushl	% ebx
movl	12(%ebp), %eax	movl	8(%ebp), %ecx
cmpl	%eax, %edx	movl	12(%ebp), %edx
jg	.16	movl	%ecx, %eax
subl	%edx, %eax	movl	%edx, %ebx
popl	% ebp	subl	%edx, %eax
ret		subl	%ecx, %ebx
.16:		cmpl	%edx, %ecx
subl	<pre>%eax, %edx</pre>	cmovle	%ebx, %eax
movl	%edx, %eax	popl	%ebx
popl	% ebp	leave	
ret		ret	

```
int distance(int a, int b)
{
   return (a > b) ? a-b : b-a;
}
```

<u>gcc -m32 -O3</u> <u>gcc -m32 -march=core2 -O3</u> <u>gcc -m64 -O3</u>

distance:

distance:

distance:

pushl	%ebp	pushl	%ebp	movl %edi, %eax
movl	%esp, %ebp	movl	%esp, %ebp	movl %esi, %edx
movl	8(%ebp), %edx	pushl	%ebx	<pre>subl %esi, %eax</pre>
movl	12(%ebp), %eax	movl	8(%ebp), %ecx	<pre>subl %edi, %edx</pre>
cmpl	%eax, %edx	movl	12(%ebp), %edx	cmpl %esi, %edi
jg	.16	movl	%ecx, %eax	<pre>cmovle %edx, %eax</pre>
subl	%edx, %eax	movl	%edx, %ebx	ret
popl	% ebp	subl	%edx, %eax	
ret		subl	%ecx, %ebx	
.16:		cmpl	%edx, %ecx	
subl	%eax, %edx	cmovle	%ebx, %eax	
movl	%edx, %eax	popl	%ebx	
popl	%ebp	leave		
ret		ret		

```
int distance(int a, int b)
{
   return (a > b) ? a-b : b-a;
}
```

<u>gcc -m32 -O3</u> <u>gcc -m32 -march=core2 -O3</u> <u>gcc -m64 -O3</u>

distance:

distance:

distance:

pushl	%ebp	pushl	%ebp	movl %edi, %eax
movl	%esp, %ebp	movl	%esp, %ebp	movl %esi, %edx
movl	8(%ebp), %edx	pushl	%ebx	<pre>subl %esi, %eax</pre>
movl	12(%ebp) , % eax	movl	8(%ebp), %ecx	<pre>subl %edi, %edx</pre>
cmpl	<pre>%eax, %edx</pre>	movl	12(%ebp), %edx	cmpl %esi, %edi
jg	.16	movl	%ecx, %eax	<pre>cmovle %edx, %eax</pre>
subl	%edx, %eax	movl	%edx, %ebx	ret
popl	%ebp	subl	%edx, %eax	
ret		subl	%ecx, %ebx	
.16:		cmpl	%edx, %ecx	
subl	<pre>%eax, %edx</pre>	cmovle	%ebx, %eax	
movl	%edx, %eax	popl	%ebx	
popl	%ebp	leave		
ret		ret		

• Note: only works (on x86) for non-memory operations

• Why? Must ensure the memory operation won't seg fault

```
int func(int a, int b, int* array)
{
  return (a > 0) ? b : array[b];
}
```

```
int func(int a, int b, int* array)
{
  return (a > 0) ? b : array[b];
}
func:
    testl %edi, %edi
    jg .L2
    movslq %esi,%rax
    movl (%rdx,%rax,4), %esi
.L2:
    movl %esi, %eax
    ret
```

- Comments on x86 assembly
 - Right-most register is usually the output register
 - "mov" instructions are loads, stores, or register copies
 - Memory access is indicated by "(...)"
 - "%exx" are 32-bit registers, "%rxx" are 64-bit registers
 - movslq: "sign extending" a "long" (32-bit) to a quad (64-bit)

```
int func2(int a, int b, int* array)
int func(int a, int b, int* array)
                                      {
                                       int temp = array[b];
  return (a > 0) ? b : array[b];
                                       return (a > 0) ? b : temp;
}
                                      }
 func:
         testl %edi, %edi
                 .L2
         jq
        movslq %esi,%rax
        movl (%rdx, %rax, 4), %esi
 .L2:
        movl
                %esi, %eax
         ret
```

- Comments on x86 assembly
 - Right-most register is usually the output register
 - "mov" instructions are loads, stores, or register copies
 - Memory access is indicated by "(...)"
 - "%exx" are 32-bit registers, "%rxx" are 64-bit registers
 - movslq: "sign extending" a "long" (32-bit) to a quad (64-bit)

```
int func2(int a, int b, int* array)
int func(int a, int b, int* array)
                                     {
                                       int temp = array[b];
  return (a > 0) ? b : array[b];
                                       return (a > 0) ? b : temp;
}
                                     }
 func:
                                     func2:
         testl %edi, %edi
                                            movslq %esi,%rax
                .L2
         jq
                                            testl %edi, %edi
        movslq %esi,%rax
                                            cmovle (%rdx,%rax,4), %esi
        movl
                (%rdx,%rax,4), %esi
                                                    %esi, %eax
                                            movl
 .L2:
                                            ret
        movl
                %esi, %eax
        ret
```

- Comments on x86 assembly
 - Right-most register is usually the output register
 - "mov" instructions are loads, stores, or register copies
 - Memory access is indicated by "(...)"
 - "%exx" are 32-bit registers, "%rxx" are 64-bit registers
 - movslq: "sign extending" a "long" (32-bit) to a quad (64-bit)

What Can the Compiler Do or Not Do?

```
void loop(int64 n, int64* array1, int64* array2, int64* ptr)
{
  for (int i = 0; i < n; i++) {
    for (int j = 0; j < n; j++) {
        if (*ptr > array2[j]) {
            array1[i] += array2[j];
        }
    }
}
```

What Can the Compiler Do or Not Do?

```
void loop(int64 n, int64* array1, int64* array2, int64* ptr)
{
   for (int i = 0; i < n; i++) {
     for (int j = 0; j < n; j++) {
        if (*ptr > array2[j]) {
            array1[i] += array2[j];
        }
     }
}
```

• How might you rewrite the code to help?

```
<u>Inner loop of loop()</u>
void loop(int64 n, int64* array1,
                                        .L5:
                                             movq (%rdx,%rax,8), %r8
          int64* array2, int64* ptr)
                                              cmpq %r8, (%rcx)
{
                                             jle .L4
  for (int64 i = 0; i < n; i++) {
                                             addq %r8, (%rsi,%r9,8)
    for (int64 j = 0; j < n; j++)  {
                                             addq $1, %rax
                                        .L4:
      if (*ptr > array2[j]) {
                                             cmpq %rax, %rdi
        array1[i] += array2[j];
                                                              load, add, store
                                              ja
                                                    .15
      }
    }
  }
```

```
void loop(int64 n, int64* array1,
                                        .L5:
          int64* array2, int64* ptr)
{
  for (int64 i = 0; i < n; i++) {
    for (int64 j = 0; j < n; j++) {
                                        .L4:
      if (*ptr > array2[j]) {
        array1[i] += array2[j];
      }
void loop2(int64 n, int64* array1,
           int64* array2, int64* ptr)
{
  for (int64 i = 0; i < n; i++) {
    int64 accumulator = 0;
    for (int 64 j = 0; j < n; j++) {
      if (*ptr > array2[j]) {
        accumulator += array2[j];
    array1[i] = accumulator;
```

Inner loop of loop()

: movq (%rdx,%rax,8), %r8
 cmpq %r8, (%rcx)
 jle .L4
 addq %r8, (%rsi,%r9,8)
: addq \$1, %rax
 cmpq %rax, %rdi
 jg .L5 load, add, store

```
<u>Inner loop of loop()</u>
void loop(int64 n, int64* array1,
                                       .L5:
                                             movq (%rdx,%rax,8), %r8
          int64* array2, int64* ptr)
                                             cmpq %r8, (%rcx)
{
                                             jle .L4
  for (int64 i = 0; i < n; i++) {
                                             addq %r8, (%rsi,%r9,8)
    for (int64 j = 0; j < n; j++) {
                                       .L4:
                                             addq $1, %rax
      if (*ptr > array2[j]) {
                                             cmpq %rax, %rdi
        array1[i] += array2[j];
                                                             load, add, store
                                                   .15
                                             ja
      }
    }
                                                Inner loop of loop2()
                                       .L13: movq (%rdx,%rax,8), %r8
                                             cmpq %r10, %r8
void loop2(int64 n, int64* array1,
                                                  .L12
                                             jge
           int64* array2, int64* ptr)
                                             addq %r8, %r9
                                       .L12: addg $1, %rax
{
  for (int64 i = 0; i < n; i++) {
                                             cmpq %rax, %rdi
    int64 accumulator = 0;
                                             jg
                                                   .L13
    for (int 64 j = 0; j < n; j++) {
      if (*ptr > array2[j]) {
        accumulator += array2[j];
    array1[i] = accumulator;
  }
```

```
Inner loop of loop()
void loop(int64 n, int64* array1,
                                       .L5:
                                            movq (%rdx,%rax,8), %r8
          int64* array2, int64* ptr)
                                             cmpq %r8, (%rcx)
{
                                             jle .L4
 for (int64 i = 0; i < n; i++) {
                                             addg %r8, (%rsi,%r9,8)
    for (int64 j = 0; j < n; j++) {
                                       .L4:
                                            addq $1, %rax
      if (*ptr > array2[j]) {
                                             cmpq %rax, %rdi
        array1[i] += array2[j];
                                                            load, add, store
                                                   .15
                                             ja
    }
                                                Inner loop of loop2()
                                       .L13: movq (%rdx,%rax,8), %r8
                                             cmpg %r10, %r8
void loop2(int64 n, int64* array1,
                                                  .L12
                                             jge
           int64* array2, int64* ptr)
                                            addq %r8, %r9
                                       .L12: addg $1, %rax
{
                                             cmpq %rax, %rdi
  for (int64 i = 0; i < n; i++) {
    int64 accumulator = 0;
                                                   .L13
                                             jg
    for (int64 j = 0; j < n; j++) {
                                            Inner loop of loop2() + CMOV
      if (*ptr > array2[j]) {
                                       .L14: movq (%r11,%r8,8), %rax
        accumulator += array2[j];
                                           ≯leaq (%r9,%rax), %rdx
                                             cmpq %r10, %rax
                                 not a load!
                                             cmovl %rdx, %r9
    array1[i] = accumulator;
                                             addq $1, %r8
                                             cmpq %r8, %rdi
                                                                  16
                                             jg
                                                   .L14
```

- To help the compiler, you can give it "profile" feedback
 - gcc -fprofile-generate ...
 - <run program>
 - gcc -fprofile-use ...

- To help the compiler, you can give it "profile" feedback
 - gcc -fprofile-generate ...
 - <run program>
 - gcc -fprofile-use ...
- Helps compiler:
 - Know which way branches commonly go
 - Frequency of loops (loop unrolling, loop peeling, etc.)

- To help the compiler, you can give it "profile" feedback
 - gcc -fprofile-generate ...
 - <run program>
 - gcc -fprofile-use ...
- Helps compiler:
 - Know which way branches commonly go
 - Frequency of loops (loop unrolling, loop peeling, etc.)
- Good just-in-time (JIT) compilers do this automatically

- Caches exploit locality
 - Temporal: same memory accessed again soon
 - **Spatial**: in same region of memory (64 byte blocks)

- Caches exploit locality
 - Temporal: same memory accessed again soon
 - **Spatial**: in same region of memory (64 byte blocks)
- Hierarchy of caches
 - First-level **instruction and data caches** (~32KB): a few cycles
 - Second/Third-level cache (~256KB to ~*MB): a dozen cycles or so
 - Main memory: ~150 cycles, highly variable
 - Limited bandwidth as well (only so many bits to/from memory)

- Caches exploit locality
 - Temporal: same memory accessed again soon
 - **Spatial**: in same region of memory (64 byte blocks)
- Hierarchy of caches
 - First-level **instruction and data caches** (~32KB): a few cycles
 - Second/Third-level cache (~256KB to ~*MB): a dozen cycles or so
 - Main memory: ~150 cycles, highly variable
 - Limited bandwidth as well (only so many bits to/from memory)
- Translation look-aside buffers (TLBs)
 - Map pages (4KB), miss can be expensive (at least dozens of cycles)

- Caches exploit locality
 - Temporal: same memory accessed again soon
 - **Spatial**: in same region of memory (64 byte blocks)
- Hierarchy of caches
 - First-level **instruction and data caches** (~32KB): a few cycles
 - Second/Third-level cache (~256KB to ~*MB): a dozen cycles or so
 - Main memory: ~150 cycles, highly variable
 - Limited bandwidth as well (only so many bits to/from memory)
- Translation look-aside buffers (TLBs)
 - Map pages (4KB), miss can be expensive (at least dozens of cycles)
- Hardware prefetching
 - Good, but makes it difficult to reason about
- **Increase locality** (temporal or spatial)
 - Change access order to re-use data
 - Pack data for spatial locality
 - Reduce code footprint (instruction cache behavior)

- **Increase locality** (temporal or spatial)
 - Change access order to re-use data
 - Pack data for spatial locality
 - Reduce code footprint (instruction cache behavior)
- Memory allocation location issues
 - Example: reduce TLB misses by packing objects onto same page

- **Increase locality** (temporal or spatial)
 - Change access order to re-use data
 - Pack data for spatial locality
 - Reduce code footprint (instruction cache behavior)
- Memory allocation location issues
 - Example: reduce TLB misses by packing objects onto same page
- Software prefetching?
 - __builtin_prefetch(const void *addr)

- **Increase locality** (temporal or spatial)
 - Change access order to re-use data
 - Pack data for spatial locality
 - Reduce code footprint (instruction cache behavior)
- Memory allocation location issues
 - Example: reduce TLB misses by packing objects onto same page
- Software prefetching?
 - __builtin_prefetch(const void *addr)

• Random versus sequential access patterns

- Pointer chasing: hard to prefetch, chain of dependencies
- Array walking: easy to prefetch, few dependencies

- **Increase locality** (temporal or spatial)
 - Change access order to re-use data
 - Pack data for spatial locality
 - Reduce code footprint (instruction cache behavior)
- Memory allocation location issues
 - Example: reduce TLB misses by packing objects onto same page
- Software prefetching?
 - __builtin_prefetch(const void *addr)

• Random versus sequential access patterns

- Pointer chasing: hard to prefetch, chain of dependencies
- Array walking: easy to prefetch, few dependencies
- Example: is a binary tree or an N-array tree more efficient?
 - Binary tree nodes are small (low spatial locality)
 - Binary trees do lots of pointer indirections

More Information

- "Intel 64 and IA-32 Architectures Optimization Reference Manual"
 - <u>http://www.intel.com/Assets/PDF/manual/248966.pdf</u>
- "Software Optimization Guide for AMD Family 10h Processors"
 - <u>http://www.amd.com/us-en/assets/content_type/</u> white_papers_and_tech_docs/40546.pdf
- "man gcc" for compiler flags
- "What Every Programmer Should Know About Memory"
 - Ulrich Drepper, <u>http://people.redhat.com/drepper/cpumemory.pdf</u>
- "Instruction latencies and throughput for x86 processors"
 - Torbjorn Granlund, <u>http://gmplib.org/~tege/x86-timing.pdf</u>
- "x86-64 Machine-Level Programming"
 - Bryant & O'Hallaron,
 - <u>http://www.cs.cmu.edu/~fp/courses/15213-s06/misc/asm64-handout.pdf</u>

Low-Level Tuning of Parallel Code

Multicore Organization

- Multiple threads per "core"
 - Share functional units, data cache, instruction cache

Multicore Organization

- Multiple threads per "core"
 - Share functional units, data cache, instruction cache
- Multiple core per chip (or "socket")
 - Private first-level caches (8KB to 64KB)
 - Maybe a second-level cache (256KB or 512KB)
 - Last-level on-chip cache (4MB to 8MB)
 - Memory controllers for directly connect memory chips (DRAM)

Multicore Organization

- Multiple threads per "core"
 - Share functional units, data cache, instruction cache
- Multiple core per chip (or "socket")
 - Private first-level caches (8KB to 64KB)
 - Maybe a second-level cache (256KB or 512KB)
 - Last-level on-chip cache (4MB to 8MB)
 - Memory controllers for directly connect memory chips (DRAM)
- Multiple sockets per system
 - Dual-socket servers standard for years
 - Quad-socket becoming common (from AMD & Intel)
 - Large systems with hundreds of sockets have been built
 - Today, usually no caching shared among sockets

• Chips connected point-to-point by high-speed links

- Chips connected point-to-point by high-speed links
- Each chip connected to 1/Nth the off-chip memory
 - Accessing a "local" page is lower latency than a "remote" page
 - Today, maybe 2x difference at most
 - Old term: Non-Uniform Memory Architecture (NUMA)
 - Was more important when latency difference was 10x or 100x

- Chips connected point-to-point by high-speed links
- Each chip connected to 1/Nth the off-chip memory
 - Accessing a "local" page is lower latency than a "remote" page
 - Today, maybe 2x difference at most
 - Old term: Non-Uniform Memory Architecture (NUMA)
 - Was more important when latency difference was 10x or 100x
- Page mapping controlled by OS
 - To optimize for lowest latency:
 - All local pages
 - To optimize for highest bandwidth:
 - Distribute pages to use all memory controllers
 - What happens in practice? oblivious allocation probably does okay

Dual-Socket Dual-Core "Core 2"



Dual-Socket Quad-Core "Core 2"



Single-Socket Quad-Core "Core i7"



26

Dual-Socket Quad-Core "Core i7"



Dual-Socket Oct-Core "Niagara T2"



- Goal: whenever reading a location in memory...
 - ...receive the value of the last write to that location

- Goal: whenever reading a location in memory...
 - ...receive the value of the last write to that location
- Invariant: Single writer --or-- one or more readers

- Goal: whenever reading a location in memory...
 - ...receive the value of the last write to that location
- Invariant: Single writer --or-- one or more readers
- Implementation
 - Each cache tracks the "state" of each cache block
 - Modified: read/write state Shared: read-only
 - Invalid: no reads or writes
 - Must transition to "modified" to be able to write block
 - Must "invalidate" all Shared copies before writing block

- Goal: whenever reading a location in memory...
 - ...receive the value of the last write to that location
- Invariant: Single writer --or-- one or more readers
- Implementation
 - Each cache tracks the "state" of each cache block
 - Modified: read/write state Shared: read-only
 - Invalid: no reads or writes
 - Must transition to "modified" to be able to write block
 - Must "invalidate" all Shared copies before writing block
- Performance implication: sharing is slow
 - Slower off-chip sharing than on-chip sharing

- Read-only sharing
 - Fine!

- Read-only sharing
 - Fine!
- True (read/write) sharing
 - Write a block invalidates all readers... so reader miss on next read
 - Next writer may stall waiting for invalidations to complete
 - Take away: avoid read/write sharing of data

- Read-only sharing
 - Fine!
- True (read/write) sharing
 - Write a block invalidates all readers... so reader miss on next read
 - Next writer may stall waiting for invalidations to complete
 - Take away: avoid read/write sharing of data
- False sharing
 - If two variables fall in the same cache block, block may ping-pong
 - Example: Each core increments "own" counter in an array
 - Repeatedly fetch and invalidate block -> **extremely slow**
 - Possible solutions:
 - Manually pad data to ensure they are in different cache blocks
 - malloc() could align each object to 64B boundaries

Multicore Memory System Implications

Multicore Memory System Implications

- Positive
 - Multithreading hides lots of miss latency!

Multicore Memory System Implications

- Positive
 - Multithreading hides lots of miss latency!
- Negative
 - Higher miss rates
 - More threads all trying to use limited shared cache capacity
 - Impact depends on data locality among threads
 - Sharing misses

• High miss latencies

- Sharing misses are often slower
- Perhaps missing to other socket's memory
- More memory bandwidth bottlenecks
 - Many parallel misses trying to use limited memory bandwidth
 - Adding memory bandwidth is expensive, so big concern



Next Time

• Review of processes and threads.

Next Time

- Review of processes and threads.
- Review of synchronization, locks and barriers
 - Shared memory/multithreaded coding
 - Synchronization primitives