Course 02158

Parallel Computation

Hans Henrik Løvengreen

DTU Compute

Ever Growing Need for Computational Power

Modelling and analysis

- Climate prediction
- System biology whole genome sequencing
- Particle collission analysis
- Market analysis

Entertainment

- Computer animation
- Virtual/augmented reality
- Speech recognition

Technology

• Machine learning

Performance Improvement Measures

Rely on Moore's Law?

- "Size and speed grows with a factor 2 every 2 years"
- Speed of a single processor has stalled at 2-3 GHz for the last 15 years
- Why? Heat!

Solution

- Many (releatively slow) processors on a single chip
- Such on-chip processing elements are often called cores
- In this course: processing element = core = processor = CPU
- We also assume processors to be *uniform* (having same capabilities)
- Recent development: *performance cores* vs. *efficiency cores* not covered

Performance

Notions

- p number of processors
- T_1 Time for sequential computation
- T_p Time for parallel computation (using p processors)

Measures

• Speedup

$$S \stackrel{\Delta}{=} \frac{T_1}{T_p}$$

• Efficiency

$$E \stackrel{\Delta}{=} \frac{S}{p} = \frac{T_1}{p \times T_p}$$



Performance Improvement Tricks

Caching

- Bulk memory is slow relative to processor speed
- Exploit *spatial* and *temporal locality* of memory accesses
- Hold *lines* of memory in fast memory the *cache*
- Many levels *memory hierarchy*

Instruction-level parallelism

- Exploit sequential nature of instruction execution
- Overlap execution of several instruction at a time *pipelining*
- Extended with branch prediction, speculative, and out-of-order execution
- Exploit the *stalling periods* of instruction execution
- Interleave execution of several instruction sequences hyper-threading
- Requires *register replication*

Parallel Programming The art of utilizing machine architectures with many parallel processors Program each processor individually? Mo — utilize OS scheduling of processes/threads Principles Data parallelism: Dividing the data into sections which can processed in parallel using the same procedure Task parallelism: Dividing the computation into phases which can be processed in parallel using dedicated procedures Approaches Manual reconstruction of program using *explicit parallelism*Decomposition into many small computation parts (tasks) — to be submitted to a *parallel computation engine*Automatic parallelization of sequential programs

OpenMP

Background

- Thread-programming on multi-processors: Much *boiler-plate code*
- 1990'ies: HW and SW vendors formed OpenMP ARB (Arch. Rev. Brd.)
- First OpenMP standard for Fortran in 1997 and C/C++ in 1998.
- Now at version 6.0 (November 2024) [964 pages!]

Idea

- To parallelize sequential programs by annotations (known as *pragmas*)
- Integrated within compilers (e.g. gcc)
- Supported by run-time library

OpenMP Directives

- In general, a *pragma* is an instruction or hint to the compiler
- C pragmas have the same form as preprocessor commands:

OpenMP Skeleton

```
#include <omp.h>
int main() {
    #pragma omp parallel
    { /* block to be executed by a thread team */
        int procs = omp_get_num_threads();
        int id = omp_get_thread_num();
        ...
    }
    /* gather result */
}
• No. of threads determined by environment variable OMP_NUM_THREADS
• May be overridden by parallel clause: ...num_threads(n)
```

```
Example: Hello OpenMP

#include <stdio.h>
#include <stdiib.h>
#include <stdlib.h>
#include <omp.h>

void Hello(void) {
    int my_rank = omp_get_thread_num();
    int thread_count = omp_get_num_threads();
    printf("Hello_from_thread_Wd_of_Wd\n", my_rank, thread_count);
}
int main(int argc, char* argv[]) {
    # pragma omp parallel
    Hello();
}
```

Scoping

- Global variables are normally *shared* in a *parallel* block
- Parallel clause ... private(x, y) introduces local copies

```
Example
```

```
int i, x, y;
x = 17;
y = 42;
# pragma omp parallel private(i, x);
for (i = 0; i < N; i++ } {
    int z ;
    x = x + 1;
    z = x + y;
  }
Here i, x, and z are private to each thread
```


Reduction

- Often results are collected in global variables protected by critical pragmas
- Parallel clause ... reduction (\oplus :x):
- \oplus is the reduction operator: +, *, &, |, &&, ||, max, min
- x is the *reduction variable*

Effect

- 1. Introduces a local copy x_i of x for each thread i
- 2. Initializes $x_i = < neutral element >$
- 3. At join: $x = x \oplus x_0 \oplus x_1 \oplus \cdots \oplus x_{p-1}$

Parallel Loop Construct

- Iterations in for loops are often fairly independent
- They may be divided into blocks for data parallelism
- Parallel loop: #pragma omp parallel for
- Next statement must be a for loop

Effect

- Creates a team of threads (as would a parallel construct)
- Distributes the iterations of the for statement among the threads
- By default, the iterations are evenly divided into large blocks distributed among the threads
- No explicit work division is needed
- Loop variable automatically becomes private
- Shared variables must still be protected or used for reduction

Example: Calculation of π

Given

$$\pi = 4 \cdot \sum_{i=0}^{\infty} (-1)^{i} \frac{1}{2i+1} = 4 \cdot \left(1 - \frac{1}{3} + \frac{1}{5} - \frac{1}{7} + \frac{1}{9} - \cdots\right)$$

Sequential program

```
double pi
double factor = -1.0;
double sum = 0.0;
int i;
for (i = 0; i< n; i++) {
  factor = -factor;
  sum += factor/(2*i + 1);
}
pi = 4.0 * sum;
```

Example: Calculation of π (ERRONEOUS) Given $\pi = 4 \cdot \sum_{i=0}^{\infty} (-1)^{i} \frac{1}{2i+1} == 4 \cdot \left(1 - \frac{1}{3} + \frac{1}{5} - \frac{1}{7} + \frac{1}{9} - \cdots\right)$ Parallel program double pi double factor = -1.0; double sum = 0.0; int i; #pragma omp parallel for reduction(+: sum) for (i = 0; i < n; i++) { factor = -factor; sum += factor/(2*i + 1); } pi = 4.0 * sum;

Example: Calculation of π (CORRECT)

Given

$$\pi = 4 \cdot \sum_{i=0}^{\infty} (-1)^{i} \frac{1}{2i+1} = = 4 \cdot \left(1 - \frac{1}{3} + \frac{1}{5} - \frac{1}{7} + \frac{1}{9} - \cdots\right)$$

Parallel program

```
double pi
double factor = -1.0;
double sum = 0.0;
int i;
#pragma omp parallel for reduction(+: sum) private(factor)
for (i = 0; i< n; i++) {
  factor = (i%2 ==0) ? 1.0 : -1.0;
  sum += factor/(2*i + 1);
}
pi = 4.0 * sum;
```


OpenMP Work Scheduling

- Default work division of parallel for not always adequate
- Loop clause ... schedule(kind[, k)] controls division

Kinds of schudulings

static	Default.	One (big)	chunk of	iterations	per thread.
--------	----------	-----------	----------	------------	-------------

- static, k Fixed round-robin distribution of k-iteration chunks
- dynamic, k Dynamic distribution of k-iteration chunks
- guided, k Proportional share of remaining iterations (at least k)
- auto System controlled
- runtime Determined by OMP_SCHEDULE environment variable

OpenMP Tasks

- Introduced in OpenMP 3.0 extended in 4.0 and 4.5
- A means for dynamic distribution of work on *irregular structures*:
 - Unbounded while loops
 - Lists and trees
 - Producer/consumer schemes
- Tasks are generated (created, spawned) from other tasks

#pragma omp task

• Execution of tasks is distributed among current team of threads

```
Example: List Processing
• Given a linked list of elements pointed to by p
• while (p != null) {
    compute(p);
    p = p->next;
  }
  #pragma omp parallel
•
    ſ
      #pragma omp single
        while (p != null) {
          #pragma omp task firstprivate(p)
            compute(p);
          p = p->next;
        }
    }
```

OpenMP Summary

- OpenMP attemps to make parallel programming *declarative*
- Still calls for *explicit* parallelism
- Supports both *explicit* and *implicit* work distribution
- Requires *compiler integraton* and a *runtime*
- Works only with C/C++ (and Fortran)
- Generalizes loop parallelism to *task-based parallelism*

Message Passing Interface (MPI) — History

$\mathsf{Background} \sim 1990$

- Many vendors of "supercomputers" with similar software
- A group of 80 people (vendors and researchers) formed MPI Forum
- Goal: A uniform way of programming *distributed memory systems*
- June 1994: Version 1.0
- API and protocols for interaction (ca. 130 functions)

Status

- De facto standard for programming distributed memory systems
- Official bindings for: C, C++, FORTRAN
- Unofficial bindings for: .NET, Java, Python, ...
- OpenMPI is a very common, open source implementation
- Current major versions: 3.1 (2015), 4.1 (2021), 5.0 (2023) ...

MPI — Key Notions

Processes

- An MPI application consists of a set of communicating (OS) processes
- Processes are (usually) *single-threaded*

— Single Program Multiple Data (SPMD)

Communicators

- A communicator is communication universe with a set of processes
- Each process within a communicator has a unique *rank* (0...)
- The full set of started processes belong to MPI_COMM_WORLD
- Processes within a communicator may communicate using:
 - Point-to-point communication
 - Collective communication

MPI — Point-to-Point Communication

Operations

- In s: MPI_Send(bufps, counts, types, dest, tags, comms)
- In r: MPI_Recv(bufpr, countr, typer, source, tagr, commr, statusp)
- Type codes: MPI_INT, MPI_CHAR, MPI_BYTE, ...
- *source* and *tag_r* may be *wildcards*: MPI_ANY_TAG, MPI_ANY_SOURCE (*)

Communication

• If operations *match*: $comm_r = comm_s, dest = r$,

$$source = s$$
 or $source = *$
 $tag_r = tag_s$, or $tag_r = *$

$$type_s = type_r$$

- Then: *count_s* elements are copied from *bufp_s* to *bufp_r*.
- If *count_s* > *count_r* an *error* occurs.
- From *statusp*, the source *s*, *tag_s*, and *count_s* may be retrieved.

MPI — Point-to-Point Semantics

General

- Communication is *reliable*
- Communication between a given sender and a given receiver is ordered
- No fairness guarantee starvation may occur

Standard mode

- MPI may buffer the message (or not).
- Both MPI_Send and MPI_Recv are blocking
- When MPI_Send returns, its buffer may be reused

Alternatives

- Other modes: Synchronous, buffered, non-blocking, ...
- Many auxiliary operations, e.g. MPI_Probe()

Example: Hello MPI World

Example: Hello MPI World

```
if (my_rank != 0) {
    sprintf(greeting, "Greetings_from_process_%d_of_%d!",
        my_rank, comm_sz);
    /* Send message to process 0 */
    MPI_Send(greeting, strlen(greeting)+1, MPI_CHAR, 0, 0,
        MPI_COMM_WORLD);
} else {
    printf("Greetings_from_process_%d_of_%d!\n", my_rank, comm_sz);
    for (int q = 1; q < comm_sz; q++) {
        /* Receive message from process q */
        MPI_Recv(greeting, MAX_STRING, MPI_CHAR, q,
        0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
    printf("%s\n", greeting);
    }
}</pre>
```

Collec	tive Commun	ication		
Basis • Barri	er:			
		:	:	:
	$MPI_Barrier(c);$	$MPI_Barrier(c);$	$MPI_Barrier(c);$	$MPI_Barrier(c);$
	•	:	:	:
• To b	e called by all in (communicator c		
Extensio	ons			
• Broa	dcast			
• Scat	tering			
• Gath	ering			
• Redu	iction			

Collective Communication — Allgather

Operation

- MPI_Allgather(*bufps*, *counts*, *types*, *bufpr*, *countr*, *typer*, *comm*)
- To be called by **all** in communicator senders/receivers
- Only *bufp*_r acts as output (for all)— all others as input
- Must all agree on $type_s = type_r$, $count_s = count_r$, comm

Effect

Collective Communication — **Gather** + **processing** = **Reduce**

Operation

- MPI_Reduce(*bufps*, *bufpr*, *count*, *type*, *op*, *dest*, *comm*)
- To be called by **all** in communicator senders and receiver
- Only *bufp*_r acts as output (for *dest*)— all others as input
- Op codes: MPI_SUM, MPI_PROD, MPI_MAX, MPI_MIN, ...
- Must all agree on type, count, op, dest, comm

Effect

