

High Performance Application Assignment Project Exam Help Optimization on CPUs

<https://powcoder.com>

Carnegie Mellon University

Add WeChat powcoder

Ian Lane

What we discussed last time:

Fast Platforms

Good Techniques

Assignment Project Exam Help

- Multicore platforms
- Manycore platforms
- Cloud platforms

- Data structures
- Algorithms
- Software Architecture

<https://powcoder.com>

Add WeChat powcoder

- Highlighted the difference between multicore and manycore platforms
- Discussed the multicore and manycore platform intricacies
- Exposing concurrency in the k-means algorithm
- Exploiting parallelism by exploring mappings from application to platform

Example 7: Sequential Optimization

```
#include <omp.h>
#include <stdio.h>
static long num_steps = 100000000; double step;
#define NUM_THREADS 2
int main ()
```

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

```
{    int i; double pi, x, local_sum; int nthreads;
    step = 1.0/(double) num_steps;
    #pragma omp parallel num_threads(NUM_THREADS) private(x, local_sum)
    {
        int i, id = omp_get_thread_num(); local_sum=0.0;
        int nthrds = omp_get_num_threads();
        double hoop = nthrds*step;
        for (i=id, x=(i+0.5)*step;i< num_steps; i=i+nthrds) {
            x += hoop;
            local_sum += 4.0/(1.0+x*x);
        }
        #pragma omp critical
        pi += local_sum * step;
    }
    printf("pi = %f \n", pi);
    return 0;
}
```

Sequential: 2.12 seconds

0.46 seconds

Outline

- Multicores Shared-Memory Platforms
- OpenMP – Parallel Shared-Memory Multithreading
 - Overview
 - Example: numerical integration of Pi
 - Shared variable synchronization
 - **SPMD vs worksharing**
 - Data environment options
 - Advanced worksharing options
 - Synchronization options
- Sequential Optimizations
 - Cost Model
 - Cost Reduction

<https://powcoder.com>

Add WeChat powcoder

SPMD vs. Worksharing

- A parallel construct by itself creates an SPMD
 - “Single Program Multiple Data”
 - Programmer must explicitly specify what each thread must do differently
 - The division of work is hard-coded in the program
- Opportunity:
 - Many parallel regions are loops
- Question:
 - Can we make it easier to parallelize these loops?

<https://powcoder.com>

Add WeChat powcoder

OpenMP: Parallel For

- One of the worksharing constructs OpenMP provide is “**omp for**”

Sequential code

```
for(i=0;i<N;i++) { a[i] = a[i] + b[i];}
```

OpenMP “parallel” region

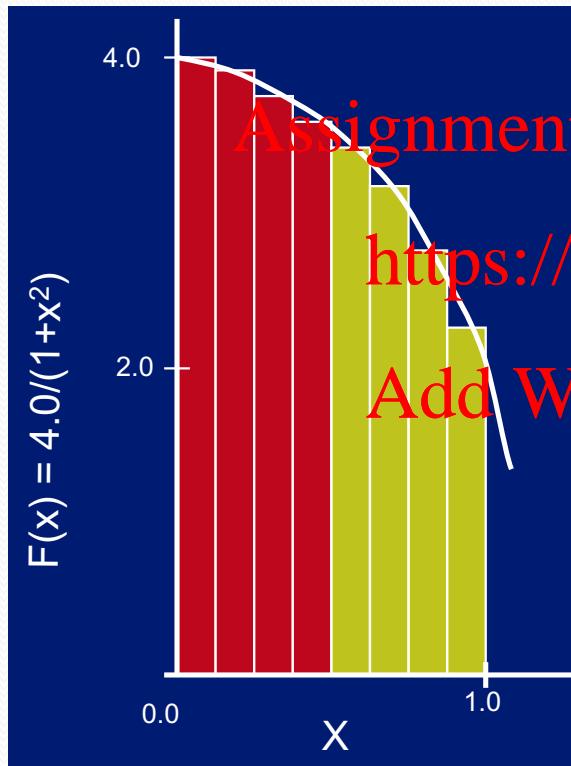
```
#pragma omp parallel
{
    int id, i, Nthrds, start, end;
    id = omp_get_thread_num();
    Nthrds = omp_get_num_threads();
    istart = id * N / Nthrds;
    iend = (id+1) * N / Nthrds;
    if (id == Nthrds-1)iend = N;
    for(i=istart;i<iend;i++) { a[i] = a[i] + b[i];}
}
```

Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder

OpenMP “parallel” region and a worksharing “**for**” construct

```
#pragma omp parallel
#pragma omp for
for(i=0;i<N;i++) { a[i] = a[i] + b[i];}
```

Other Concurrency Opportunities

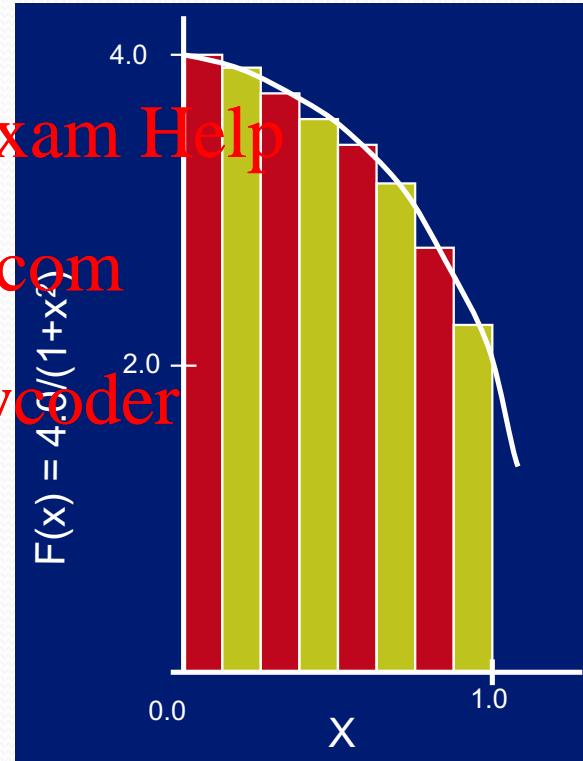


Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

- █ Thread 0
- █ Thread 1



Example 8: Parallel For and Reduction

```
#include <omp.h>
#include <stdio.h>
static long num_steps = 100000000; double step;
#define NUM_THREADS 2
int main ()
```

[Assignment Project Exam Help](https://powcoder.com)
<https://powcoder.com>
Add WeChat powcoder

```
{    int i; double pi, x, sum = 0.0;
    step = 1.0/(double) num_steps;
    #pragma omp parallel private(x)
    {
        #pragma omp for reduction(+:sum)
        for (i=0;i< num_steps;i++)
            x = (i+0.5)*step;
            sum = sum + 4.0/(1.0+x*x);
    }
    Pi = sum * step;
    printf("pi = %f \n", pi);
    return 0;
}
```

Sequential: 2.12 seconds

1.22 seconds

Reduction is more scalable than critical sections.

“i” is private by default

Loops must not have loop carry dependency.
Limitation of compiler technology.

Combined Parallel/Worksharing Construct

- OpenMP shortcut: Put the “parallel” and the worksharing directive on the same line

Assignment Project Exam Help

```
double res[MAX]; int i;  
#pragma omp parallel  
{  
    #pragma omp for  
    for (i=0;i< MAX; i++) {  
        res[i] = huge();  
    }  
}
```

<https://powcoder.com>

Add WeChat powcoder



```
double res[MAX]; int i;  
#pragma omp parallel for  
for (i=0;i< MAX; i++) {  
    res[i] = huge();  
}
```

OpenMP: Working With Loops

- Basic approach
 - Find compute intensive loops
 - Make the loop iterations independent. So they can safely execute in any order without loop-carried dependencies
 - Place the appropriate OpenMP directive and test

<https://powcoder.com>

```
double hoop = nthrds*step;
for (i=id, x=(i+0.5)*step;i< num_steps; i=i+nthrds) {
    x += hoop;
    sum = sum + 4.0/(1.0+x*x);
}
```

*x depend on
i not in for*



```
for (i=0;i< num_steps; i++){
    x = (i+0.5)*step;
    sum = sum + 4.0/(1.0+x*x);
}
```

OpenMP: Reductions

```
for (i=0;i< num_steps; i++){
    x = (i+0.5)*step;
    sum = sum + 4.0/(1.0+x*x);
}
```

Assignment Project Exam Help



```
#pragma omp for reduction(+:sum)
for (i=0;i< num_steps; i++){
    x = (i+0.5)*step;
    sum = sum + 4.0/(1.0+x*x);
}
```

<https://powcoder.com>

Add WeChat powcoder

reduction (op : list)

1. A local copy of each list variable is made and initialized depending on the “op” (e.g. 0 for “+”)
2. Updates occur on the local copy
3. Local copies are reduced into a single value and combined with the original global value

- Accumulating values into a single variable (sum) creates true dependence between loop iterations that can't be trivially removed
- This is a very common situation ... it is called a “reduction”.
- Support for reduction operations is included in most parallel programming environments.

OpenMP: Reduction Operators

- Many different associative operands can be used with reduction:
- Initial values are the ones that make sense mathematically.

Assignment Project Exam Help

| Operator | Initial value |
|----------|---------------|
| + | 0 |
| * | 1 |
| - | 0 |

| Operator | Initial value |
|----------|---------------|
| & | ~ 0 |
| | 0 |
| ^ | 0 |
| && | 1 |
| | 0 |

Example 8: Parallel For and Reduction

```
#include <omp.h>
#include <stdio.h>
static long num_steps = 100000000; double step;

int main () {
    int i; double pi, x, sum = 0.0;
    step = 1.0/(double) num_steps;
    #pragma omp parallel for private(x) reduction(+:sum)
    for (i=0;i< num_steps; i++){
        x = (i+0.5)*step;
        sum = sum + 4.0/(1.0+x*x);
    }
    Pi = sum * step;
    printf("pi = %f \n", pi);
    return 0;
}
```

Sequential: 2.12 seconds



1.22 seconds

<https://powcoder.com>

Add WeChat powcoder

Advantage: we created a parallel program without changing any code and by adding 2 simple lines of text!

Disadvantage:

- 1) Lot's of pitfalls if you don't understand system architecture
- 2) The basic result may not be the fastest one can do

Outline

- Multicores Shared-Memory Platforms
- OpenMP – Parallel Shared-Memory Multiprocessing
 - Overview
 - Example: numerical integration of Pi
 - Shared variable synchronization
 - SPMD vs worksharing
 - **Data environment options**
 - Advanced worksharing options
 - Synchronization options
- Sequential Optimizations
 - Cost Model
 - Cost Reduction

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Data Environment in OpenMP

- Shared Memory programming model:
 - Most variables are shared by default

Assignment Project Exam Help

- Global variables are **SHARED** among threads
 - File scope variables, static
 - Dynamically allocated memory (ALLOCATE, malloc, new)

Add WeChat powcoder

- But not everything is shared...
 - Functions called from parallel regions are **PRIVATE**
 - Automatic variables within a statement block are **PRIVATE**.

Example: Data Sharing

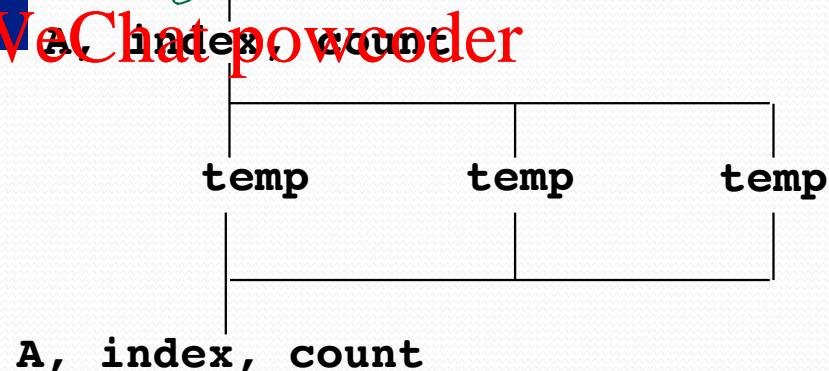
```
double A[10];
int main() {
    int index[10];
    #pragma omp parallel
        work(index);
    printf("%d\n", index[0]);
}
```

Assignment Project Exam Help
<https://powcoder.com>

extern double A[10];
void work(int *index) {
 double temp[10];
 static int count;
 ...
}

A, index and count are shared by all threads.

temp is local to each thread



Changing Storage Attributes

- One can selectively change storage attributes for constructs using the following clauses:

- SHARED
- PRIVATE
- FIRSTPRIVATE

Assignment Project Exam Help

<https://powcoder.com>

- The final value of a private inside a parallel loop can be transmitted to the shared variable outside the loop with:

- LASTPRIVATE

Add WeChat powcoder

typically use critical / atom(...)

Example 9: firstprivate example

```
#include <omp.h>
#include <stdio.h>
static long num_steps = 100000000; double step;
#define NUM_THREADS 2
int main ()
{
    int i; double pi, x, local_sum = 0.0; int nthreads;
    step = 1.0/(double) num_steps;
    omp_set_num_threads(NUM_THREADS);
    #pragma omp parallel private(x) firstprivate(local_sum)
    {
        int i, id = omp_get_thread_num(); local_sum=0.0;
        int nthrds = omp_get_num_threads();
        double hoop = nthrds*step;
        for (i=id, x=(i+0.5)*step; i< num_steps; i=i+nthrds) {
            x += hoop;
            local_sum += 4.0/(1.0+x*x);
        }
        #pragma omp critical
        pi += local_sum * step;
    }
    printf("pi = %f \n", pi);
    return 0;
}
```

Outline

- Multicores Shared-Memory Platforms
- OpenMP – Parallel Programming Model for Shared-Memory Multiprocessors
Assignment Project Exam Help
 - Overview
 - Example: numerical integration of Pi
 - Shared variable synchronization
 - SPMD vs worksharing
 - Data environment options
 - **Advanced worksharing options**
 - Synchronization options
- Sequential Optimizations
 - Cost Model
 - Cost Reduction

<https://powcoder.com>

Add WeChat powcoder

The Scheduling Cause

- The schedule clause affects how loop iterations are mapped onto threads

- schedule(static,[chunk])
 - Deal-out blocks of iterations of size “chunk” to each thread.
- schedule(dynamic,[chunk])
 - Each thread grabs “chunk” iterations off a queue until all iterations have been handled.
- schedule(guided,[chunk])
 - Threads dynamically grab blocks of iterations. The size of the block starts large and shrinks down to size “chunk” as the calculation proceeds.
- schedule(runtime)
 - Schedule and chunk size taken from the OMP_SCHEDULE environment variable (or the runtime library)

The Scheduling Cause

| Schedule Clause | When To Use | Help |
|-----------------|---|--|
| STATIC | Pre-determined and predictable by the programmer https://powcoder.com | Least work at runtime : scheduling done at compile-time |
| DYNAMIC | Unpredictable, highly variable work per iteration Add WeChat powcoder | Most work at runtime : complex scheduling logic used at run-time |
| GUIDED | Special case of dynamic to reduce scheduling overhead | |

Example 9: Scheduling

```
#include <omp.h>
#include <stdio.h>
static long num_steps = 100000000; double step;

int main () {
    int i; double pi, x, sum = 0.0;
    step = 1.0/(double) num_steps;
    #pragma omp parallel for private(x) reduction(+:sum) schedule(static,
    10000000)
    for (i=0;i< num_steps; i++){
        x = (i+0.5)*step;
        sum = sum + 4.0/(1.0+x*x);
    }
    Pi = sum * step;
    printf("pi = %f \n", pi);
    return 0;
}
```

Sequential: 2.12 seconds

1.31 seconds

Assignment Project Exam Help

<https://powcoder.com>

#pragma omp parallel for private(x) reduction(+:sum) schedule(static,

10000000)

for (i=0;i< num_steps; i++) {

 x = (i+0.5)*step;

 sum = sum + 4.0/(1.0+x*x);

}

Pi = sum * step;

printf("pi = %f \n", pi);

return 0;

}

10 chunks -

Add WeChat powcoder

Outline

- Multicores Shared-Memory Platforms
- OpenMP – Parallel Programming Model for Shared-Memory Multiprocessors
Assignment Project Exam Help
 - Overview
 - Example: numerical integration of Pi
 - Shared variable synchronization
 - SPMD vs worksharing
 - Data environment options
 - Advanced worksharing options
 - **Synchronization options**
- Sequential Optimizations
 - Cost Model
 - Cost Reduction

<https://powcoder.com>

Add WeChat powcoder

Synchronization: ordered

- The **ordered** region executes in the sequential order
- Important for some scientific code and optimization code

Order of reduction may cause rounding differences

<https://powcoder.com>

```
#pragma omp parallel private (tmp)
#pragma omp for ordered reduction(+:res)
for (l=0;l<N;l++){
    tmp = NEAT_STUFF(l);
    #pragma ordered
    res += consum(tmp);
}
```

Add WeChat powcoder

Synchronization Barrier

- **Barrier:** Each thread waits until all threads arrive

Assignment Project Exam Help

```
#pragma omp parallel shared (A, B, C) private(id)
{
    id=omp_get_thread_num();
    A[id] = big_calc1(id);
#pragma omp barrier
#pragma omp for Add WeChat powcoder
    for(i=0;i<N;i++){C[i]=big_calc3(i,A);}
#pragma omp for nowait
    for(i=0;i<N;i++){ B[i]=big_calc2(C, i); }
    A[id] = big_calc4(id);
}
```

implicit barrier at the end of
a for worksharing construct

No implicit barrier due to
nowait

implicit barrier at the end of
a parallel region

“Single” Worksharing Construct

- The **single** construct denotes a block of code that is executed by only one thread (not necessarily the master thread).
- A barrier is implied at the end of the single block
 - can remove the barrier with a nowait clause

Assignment Project Exam Help

<https://powcoder.com>

```
#pragma omp parallel  
{  
    do_many_things();  
    #pragma omp single  
    {  
        exchange_boundaries();  
    }  
    do_many_other_things();  
}
```

Add WeChat powcoder

```
#pragma omp parallel  
do_many_things();  
#pragma omp single nowait  
{  
    exchange_boundaries();  
} do_many_other_things();  
}
```

Nested Parallelism

- **Q:** Is nested parallel possible with OpenMP?
- **A:** Yes. But be sure to understand why you want to use it.

```
#include <omp.h>
#include <stdio.h>
void report_num_threads(int level) {
    #pragma omp single
    { printf("L %d: # threads in team %d\n", level, omp_get_num_threads()); }
}
int main() {
    omp_set_dynamic(0);
    #pragma omp parallel num_threads(2)
    {
        report_num_threads(1);
        #pragma omp parallel num_threads(2)
        {
            report_num_threads(2);
            #pragma omp parallel num_threads(2)
            report_num_threads(3);
        }
    }
    return(0);
}
```

Add WeChat powcoder

<http://download.oracle.com/docs/cd/E19205-01/819-5270/aewbc/index.html>

Nested Parallelism

- Compiling and running this program with nested parallelism enabled produces the following (sorted) output:

Assignment Project Exam Help

```
$ export OMP_NESTED=TRUE  
$ ./experimentN  
L 1: # threads in team 2  
L 2: # threads in team 2  
L 2: # threads in team 2  
L 3: # threads in team 2
```

```
$ export OMP_NESTED=FALSE  
$ ./experimentN  
L 1: # threads in team 2  
L 2: # threads in team 1  
L 3: # threads in team 1  
L 2: # threads in team 1  
L 3: # threads in team 1
```

<https://powcoder.com>

Add WeChat powcoder

OpenMP Environment Variables

- OMP_SCHEDULE=algorithm
 - dynamic[, n]
 - guided[, n]
 - Runtime
 - static[, n]
- OMP_NUM_THREADS=num
- OMP_NESTED=TRUE|FALSE
- OMP_DYNAMIC=TRUE|FALSE

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Outline

- Multicores Shared-Memory Platforms
- OpenMP – Parallel Shared-Memory Multiprocessors
 - Overview
 - Example: numerical integration of Pi
 - Shared variable synchronization
 - SPMD vs worksharing
 - Data environment options
 - Advanced worksharing options
 - Synchronization options
- Sequential Optimizations
 - Cost Model
 - Cost Reduction

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Optimizing the Base Case

- Base case is usually the guts of the algorithm that computes on data in caches

Assignment Project Exam Help

- Should fit in registers

<https://powcoder.com>

- Must use SIMD

Add WeChat powcoder

Cost Model

- **Algebraic model**

- +, -, /, *

- sin, cos, etc.

```
int func(int i) {  
    return i+2;  
}
```

- **Realistic model**

- +, -, /, *
- std libs (sin, cos, etc)
- array[i]
- A.sum, A.prod
- (double) x
- func(...)
- i < n, compute
- i < n, test and branch
- {} unconditional branch

```
struct A {  
    double sum, prod;  
};  
A acc(int n, int array[]) {  
    A res = {0.0, 1.0};  
    for(int i=0; i<n; ++i) {  
        A.sum = A.sum +  
            (double)func(array[i]);  
        A.prod = A.prod *  
            (double)func(array[i]);  
    }  
    return res;  
}
```

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Cost Model (cont'd)

- Algebraic model

cost = **3 n**

```
int func(int i) {  
    return i+2;  
}
```

- Realistic model

+ and *
2 int -> double

4 a.x

++i

i<n conditional jump

2 func() calls

cost = **19 n**

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

```
struct A {  
    double sum, prod;  
};  
  
A acc(int n) {  
    A res = {0.0, 1.0};  
    for(int i=0; i<n; ++i) {  
        A.sum = A.sum + (double)func(i);  
        A.prod = A.prod * (double)func(i);  
    }  
    return res;  
}
```

6x discrepancy in cost

Using the Cost Model to Optimize

```
int func(int i) { return i+2; }
```

```
struct A { double sum, prod; };
```

Assignment Project Exam Help

```
A acc(int n, int array[]) {
```

```
    A res = {0.0, 1.0};
```

<https://powcoder.com>

```
    for(int i=0; i<n; ++i) {
```

```
        A.sum = A.sum + (double)func(i);
```

```
        A.prod = A.prod * (double)func(i);
```

```
    }
```

```
    return res;
```

```
}
```

Cost = 19n

Add WeChat powcoder

Using the Cost Model to Optimize I

```
int func(int i) { return i+2; }
```

```
struct A { double sum, prod; };
```

Assignment Project Exam Help

```
A acc(int n, int array[]) {
```

```
    double sum=0.0, prod=1.0;
```

<https://powcoder.com>

```
    for(int i=0; i<n; ++i) {
```

```
        f = array[i];
```

```
        sum = sum + f;
```

```
        prod = prod * f;
```

```
    }
```

```
    A res = {sum, prod}
```

```
    return res;
```

```
}
```

Cost = $6n$

Add WeChat powcoder

Cost Reduction = closing down the gap

Using the Cost Model to Optimize II

```
int func(int i) { return i+2; }
```

```
struct A { double sum, prod; };
```

Assignment Project Exam Help

```
A acc(int n, int array[]) {
```

```
    double sum=0.0, prod=1.0;
```

<https://powcoder.com>

```
    for(int i=2; i<n+2; ++i) {
```

```
        f = (double)(1)
```

```
        sum = sum + f;
```

```
        prod = prod * f;
```

```
    }
```

```
    A res = {sum, prod};
```

```
    return res;
```

```
}
```

Cost = 5n

Add WeChat powcoder

Optimizations enable new optimizations

Using the Cost Model to Optimize III

```
int func(int i) { return i+2; }
```

```
struct A { double sum, prod; };
```

Assignment Project Exam Help

```
A acc(int n, int array[]) {
    double sum=0.0, prod=1.0;
    #pragma unroll(4)
    for(int i=2; i<n+2; ++i) {
        f = (double)(i)
        sum = sum + f;
        prod = prod * f;
    }
    A res = {sum, prod};
    return res;
}
```

Cost = $3.5n$

Add WeChat powcoder

https://powcoder.com

Final result: $3.5n$ vs $19n$

Outline

- Multicores Shared-Memory Platforms
- OpenMP – Parallel Shared-Memory Multiprocessing
 - Overview
 - Example: numerical integration of Pi
 - Shared variable synchronization
 - SPMD vs worksharing
 - Data environment options
 - Advanced worksharing options
 - Synchronization options
- Sequential Optimizations
 - Cost Model
 - **Cost Reduction**

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Cost Reduction: A Compiler Problem

- **Solution:** “Human Compiler”
- Compilers often fail on complex codes (many assumptions are violated)

Assignment Project Exam Help

- Optimizations
 - Strength reduction
 - Function inlining
 - Loop unrolling
 - Common subexpression elimination
 - Load/store elimination
 - Table lookups
 - Branch elimination

Add WeChat powcoder

Strength Reduction

```
for i = 1..n
```

```
    sum = sum + a[i] / c;
```

Assignment Project Exam Help

<https://powcoder.com>

double c_inv = 1/c;
for i = 1..n
 sum = sum + a[i] * c_inv;

division is
more expensive
than multiplication

Expensive operations:
/, %, sin, cos, log

Function Inlining

```
double f(a) { return a/c; }
```

Assignment Project Exam Help

```
for i = 1..n
```

```
    sum = sum + f(a[i])
```

<https://powcoder.com>

Add WeChat powcoder

```
for i = 1..n
```

```
    sum = sum + a[i] / c;
```

Now strength reduction can apply

Loop Unrolling

```
for i = 1..n
```

```
    sum = sum + a[i]
```

Assignment Project Exam Help

<https://powcoder.com>

for i = 1:4:n
 sum = sum + a[i] + a[i+1] + a[i+2] + a[i+3]

```
// remaining iterations
for i = i:n
    sum = sum + a[i]
```

sequential locality

Often enables other optimizations

Common Sub-expression Elimination

```
for j = 1..m
```

```
for i = 1..n
```

```
sum = sum + cos(j*PI/180)*a[i];
```

Assignment Project Exam Help

<https://powcoder.com>



```
for j = 1..m
```

```
double c = cos(j*PI/180);
```

```
for i = 1..n
```

```
sum = sum + c*a[i];
```

Add WeChat powcoder

Table Lookups

```
for j = 1..m
```

```
    for i = 1..n
```

```
        sum = sum + cos(j*PI/180)*a[i];
```

Assignment Project Exam Help

<https://powcoder.com>



Add WeChat powcoder

```
for j = 1..m
```

```
    double c = COS_TAB[j];
```

```
    for i = 1..n
```

```
        sum = sum + c*a[i];
```

Load/Store Elimination

Loop Merging

Assignment Project Exam Help

```
for i = 1..n
    sum = sum + f(a[i]);
```

```
for i = 1..n
    prod = prod * f(a[i]);
```

<https://powcoder.com>

Add WeChat powcoder



```
for i = 1..n
    sum = sum + f(a[i]);
    prod = prod * f(a[i]);
```

One of the most important optimizations!
Almost always enables other optimizations.

Branch Elimination

```
for i = 1..n  
if(i != except)  
    sum = sum + a[i];
```

Assignment Project Exam Help

<https://powcoder.com>

```
for i = 1..except-1  
    sum = sum + a[i];  
for i = except+1..n  
    sum = sum + a[i];
```

a[except] = 0.0;
for i = n
 sum = sum + a[i];

What we have discussed in the last 2 lectures

Fast Platforms

Assignment Project Exam Help

Good Techniques

- Multicore platforms
- Manycore platforms
- Cloud platforms

- Data structures
- Algorithms
- Software Architecture

<https://powcoder.com>

- OpenMP Abstractions:
 - Help establish a mental model for the programmers
 - Make it easier to write parallel code
 - Performance depend on deep understanding of the implementation platform