CS420 – Lecture 7

Marc Snir

Fall 2018

] [

- IV	larc	-51	٦IT

≣ ► ≣ ৩৭ে Fall 2018 1/47

▲ロト ▲圖ト ▲国ト ▲国ト

Examples

■ ■ つへで Fall 2018 2/47

▲□▶ ▲圖▶ ▲厘▶ ▲厘▶

Matrix Product

- Can allocate to each processor a tile to compute provided the computation of distinct tiles are independent
- If tile in k dimension need to add a reduction.
- Always want to parallelize outermost loop (get large tasks)





If tile in k dimension, need to add reduction, $\sum_{k=1}^{n}$

```
#include <omp.h>
#include <stdio.h>
#define N 500
double a[N][N],b[N][N],c[N][N];
int main(int argc, char *argv[]) {
   int i, j, k, n;
   double time;
   for (n=0; n<10; n++) {
     time = omp_get_wtime();
     omp_set_num_threads(4);
</pre>
```

- omp_get_wtime returns time in seconds.
- Only outermost (i) loop is executed in parallel

}

▲□▶ ▲圖▶ ▲ 圖▶ ▲ 圖▶ ― 圖 … のへで

```
Tile j
```

```
#pragma omp parallel for
for (j=0; j<N; i++)
for (i=0; i<N; j++)
for (k=0; k<N; k++)
a[i][j] += b[i][k] * c[k][j];
```

Tile k

```
• The reduction clause takes an array
section argument (will be discussed later)
Tile both i and j
```

```
#pragma omp parallel for collapse(2)
for (i=0; i<N; i++)
  for (j=0; j<N; j++)
    for (k=0; k<N; k++)
        a[i][j] += b[i][k] * c[k][j];</pre>
```

 collapse(2) indicates that two outermost loops should be taken as one loop (with N×N iterates) and executed in parallel

イロト イヨト イヨト イヨト



Tile i: Each thread computes product of horizontal tile of b with c that yields a horizontal tile of a.

< □ > < 同 >

< ∃ > < 3



Tile j: Each thread computes product of b with vertical tile of c that yields a vertical tile of a.



Tile i, j: Each thread computes product of horizontal tile of b with a vertical slice of c that yields a 2D tile of a.

Fall 2018 8 / 47

< □ > < @ >

★ ∃ ►



Tile k: Each thread computes product of vertical tile of b with a horizontal slice of c that yields an $N \times N$ matrix; the resulting matrices need to be added.

Code	Time (msec)	
Tile i	1277±57	
Tile j	1623 ± 43	
Tile i,j	992±14	
Tile k	bus error	

- Reduction on array sections is new feature may not be well-supported
- Tiling choices impact locality

イロト イヨト イヨト イヨト

Sequential Jacobi

<ロト < 部 ト < 目 ト 4 目 ト 目 の Q () Fall 2018 11 / 47

Parallel Jacobi

collapse(2): the two outer loops are handled as one parallel loop with MN iterations

Possibly better

Tile the nested loops and allocate to threads full tiles

stencil computation



Tiling



Marc Snir

Fall 2018 13 / 47

★ ∃ ►

```
. . .
do {
  err = 0:
  k = 1 - k:
  #pragma omp parallel for reduction(max:err)
  for (jj=1; jj<N-1; jj += T)
    for (i=1; i<M-1; i++)</pre>
      for (j=jj; j<jj+T; j++) {</pre>
        a[1-k][i][j] = 0.25 * (a[k][i-1][j] + a[k][i+1][j] +
                                  a[k][i][j-1] + a[k][i][j+1]);
        err = fmax(err, fabs(a[1][i][j]-a[0][i][j]));
      }
} while (err > maxerr);
. . .
```

Only outer loop is executed in parallel

Tiling provides the same improvements in cache hit ratio as for sequential code Assuming tiles are cache line aligned

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

```
Solution 1: Use 2D tiles
\* N=n*T1+2, M=m*T2+2 *\
. . .
do {
  err = 0; k = 1-k;
  #pragma omp parallel for \setminus
                collapse(2) reduction(max:err)
  for (ii=1; ii<N-1; ii += T1)</pre>
    for (jj=1; jj<M-1; jj += T2)</pre>
      for (i=ii; i<ii+T2; i++)</pre>
         for (j=jj; j<jj+T2; j++) {</pre>
           a[1-k][i][j] = 0.25 * (a[k][i-1][j] + a[k][i+1][j] +
                                     a[k][i][j-1] + a[k][i][j+1]);
           err = fmax(err, fabs(a[1][i][j]-a[0][i][j]));
         3
} while (err > maxerr);
. . .
                                                                イロト イヨト イヨト イヨト
```

2D Tiling

∃ nar

Solution 2: Use nested parallelism

```
. . .
do {
  err = 0:
  k = 1 - k;
  #pragma omp parallel for reduction(max:err)
  for (jj=0; jj<n; jj++) {</pre>
    #pragma omp parallel for reduction(max:err)
    for (i=0; i<M; i++) {</pre>
      for (j=jj*T1+1; j<(jj+1)*T1+1; j++) {</pre>
        a[1-k][i][j] = 0.25 * (a[k][i-1][j] + a[k][i+1][j] +
                                  a[k][i][j-1] + a[k][i][j+1]);
         err = fmax(err, fabs(a[1][i][j]-a[0][i][j]));
      }
    }
  3
} while (err > maxerr);
. . .
```

< 日 > < 同 > < 三 > < 三 >

- Might not be supported (get error) controlled by ICV
 - OMP_MAX_ACTIVE_LEVELS, omp_set_max_active_levels, omp_get_max_active_levels
- Even if it is supported it id not obvious how many threads will execute a nested loop could be one

< ロ > < 同 > < 回 > < 回 > < 回 >

Let's experiment

```
void report_num_threads(int level) {
  printf("Level %d: team size = %d n",
         level, omp_get_num_threads());
}
int main(int argc, char *argv[]) {
  printf("available threads = %d\n",
         omp_get_max_threads());
 #pragma omp parallel num_threads(2)
  Ł
    report_num_threads(1);
    #pragma omp parallel num_threads(2)
    ſ
      report_num_threads(2);
      #pragma omp parallel num_threads(2)
      report_num_threads(3);
    }
```

Executed:

export OMP_MAX_THREADS=7
./a.out

Output was:

available threads = 7 Level 1: team size = 2 Level 1: team size = 2 Level 2: team size = 1 Level 3: team size = 1 Level 2: team size = 1 Level 3: team size = 1

イロト イヨト イヨト イヨト

} } Use it at your own peril: Many implementations do not use "spare" threads to increase parallelism at lower levels Note:

```
#pragma omp parallel for
is equivalent to
#pragma omp parallel
#pragma omp for
```

first statement creates team; second statement does work sharing across team Problem: Nested team creation might use only parent thread

イロト イ理ト イヨト イヨト

SparseMV: a = b + Cd where C is a sparse matrix: most entries are zero.

Marc Snir

Fall 2018 20 / 47

- How does one store the matrix so that only non-zeros are stored?
- How does on avoid the multiplications by zero?

CRS: Compressed Row Storage

3

< □ > < □ > < □ > < □ > < □ >

SparseMV

```
for (i=0; i<N; i++)
a[i] = b[i];
for (j=row_ptr[i]; j<row_ptr[i+1]; j++)
a[i] += val[j] * d[col_idx[j]]</pre>
```


If rows have roughly the same number of non-zeros, then get good load balancing by statically tiling rows

```
#pragma omp parallel for schedule(static,T)
for (i=0; i<N; i++) {
    a[i] = b[i];
    for (j=row_ptr[i]; j<row_ptr[i+1]; j++)
        a[i] += val[j] * d[col_idx[j]];
}</pre>
```

T chosen so that tasks are large enough

< ロ > < 同 > < 回 > < 回 >

If rows have very different number of non-zeros, need to use dynamic load balancing.

```
#pragma omp parallel for schedule(dynamic,T)
for (i=0; i<N; i++) {
    a[i] = b[i];
    for (j=row_ptr[i]; j<row_ptr[i+1]-1; j++)
        a[i] += val[j] * d[col_idx[j]];
}</pre>
```

T chosen so that tasks are large enough, but number of tasks still large wrt number of threads

Gauss-Seidel

Like Jacobi, except done in place (one array)

$$a_{i,j}^{(k+1)} = 0.25(a_{i-1,j}^{(k+1)} + a_{i,j-1}^{(k+1)} + a_{i+1,j}^{(k)} + a_{i,j+1}^{(k)})$$

Sequential code

How do we parallelize?

How can we reorder the nested loop so as to have many independent operations?

イロト イ団ト イヨト イヨト 二日

Wavefront

Loop carried dependencies

Wavefront

Loop carried dependencies

Wavefront

メロト メポト メヨト メヨト

code


```
. . .
/* number of diagonals is M+N-5 */
for (d=0; d<M+N-5; d++) {
  /* first & last diagonal row */
  ifirst = (d<M-1) ? d+1 : M-2;
  ilast = (d < N - 1) ? 1 : d - M + 3;
  #pragma omp parallel for
  for (i=ifirst; i<=ilast; i++) {</pre>
    j = d+2-i;
    a[i][j] = 0.2 * (a[i][j])
                                 +
         a[i-1][j] + a[i+1][j] +
         a[i][j-1] + a[i][j+1]);
  }
```

. . .

Fall 2018 27 / 47

イロト イヨト イヨト イヨト

- Specify the set of iterations that need to be executed
- Specify the dependencies that need to be obeyed.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

```
/* both loops collapsed, must obey ordering constraints */
#pragma omp for collapse(2) ordered(2)
for (i=1; i<N-1; i++)</pre>
  for (j=1; j<M-1; j++) {</pre>
    /* must wait until (i-1,j) and (i,j-1) iterations complete */
    #pragma omp ordered depend(sink: i-1,j) depend(sink: i,j-1)
    a[i][j] = 0.2 * (a[i-1][j] + a[i+1][j] +
                     a[i][j-1] + a[i][j+1] + a[i][j]);
    /* iteration (i,j) complete, let dependencies proceed */
    #pragma omp ordered depend(source)
  }
```

- Must likely inefficient because dependencies tracked at fine grain
- Can tile

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

<ロト < 部ト < 語ト < 語ト 語 の Q () Fall 2018 30 / 47

Back to Jacobi

. . .

```
do {
  err = 0;
  k = 1 - k;
  #pragma omp parallel for reduction(max:err)
  for (jj=1; jj<N-1; jj += T)</pre>
    for (i=1; i<M-1; i++)</pre>
      for (j=jj; j<jj+T; j++) {</pre>
        a[1-k][i][j] = 0.25 * (a[k][i-1][j] + a[k][i+1][j] +
                                  a[k][i][j-1] + a[k][i][j+1]);
        err = fmax(err, fabs(a[1][i][j]-a[0][i][j]));
      }
} while (err > maxerr);
. . .
```


イロト イヨト イヨト イヨト

= 990

• Do we have races?

< ■ ト ■ シへぐ Fall 2018 32 / 47

▲□▶ ▲圖▶ ▲厘▶ ▲厘▶

- Do we have races?
- No
 - err is a reduction variable
 - During each iteration of the while loop we read one copy of a and write another copy no conflicts
 - No thread starts next iteration of the while loop before all threads completed the previous iteration there is an implicit barrier at the end of the parallel section

• Do we have communication between threads?

▲ ■ ▶ ■ シへの Fall 2018 33 / 47

- Do we have communication between threads?
- Yes "orange" thread needs a column written by each of "purple" and "blue" threads at previous "while" iteration, assuming threads pick same chunk at successive iterations.

- Do we have communication between threads?
- Yes "orange" thread needs a column written by each of "purple" and "blue" threads at previous "while" iteration, assuming threads pick same chunk at successive iterations.
- Are we sure that a thread picks same slice at successive iteration?

- Do we have communication between threads?
- Yes "orange" thread needs a column written by each of "purple" and "blue" threads at previous "while" iteration, assuming threads pick same chunk at successive iterations.
- Are we sure that a thread picks same slice at successive iteration?
- Not in general: allocation may change from parallel loop to parallel loop

Allocation does not change from one parallel for to the next if

- Number of threads is fixed
- Schedule is static

```
\times N=n*T+2 *
. . .
omp_set_dynamic(0);
do {
  err = 0:
  k = 1 - k;
  #pragma omp parallel for schedule(static) reduction(max:err)
  for (jj=1; jj<N-1; jj += T)</pre>
    for (i=1; i<M-1; i++)</pre>
      for (j=jj; j<jj+T; j++) {</pre>
         a[1-k][i][j] = 0.25 * (a[k][i-1][j] + a[k][i+1][j] +
                                  a[k][i][j-1] + a[k][i][j+1]);
         err = fmax(err, fabs(a[1][i][j]-a[0][i][j]));
      }
} while (err > maxerr);
. . .
```

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Jacobi "static" style

Can we avoid the overhead of repeatedly forking and joining control? #pragma omp parallel Ł n = omp_get_num_threads(); myid = omp_get_thread_num(); myfirst = myid*N/n; nextfirst = (myid+1)*N/n; do { #pragma omp single err = 0: mverr = 0; k = 1-k;for (i=1; i<M-1; i++)</pre> for (j=myfirst; j<nextfirst; j++) {</pre> a[1-k][i][j] = 0.25 * (a[k][i-1][j] + a[k][i+1][j] + a[k][i][j-1] + a[k][i][j+1]); myerr = fmax(myerr, fabs(a[1][i][j]-a[0][i][j])); } #pragma omp critical err = fmax(err, myerr); #pragma omp barrier } while (err > maxerr);

}

▲□▶ ▲圖▶ ▲ 圖▶ ▲ 圖▶ ― 圖 … のへで

- Code has become more verbose.
- No load balancing by OpenMP runtime user has to do it, if needed
- Cannot use reduction
- \bullet Cannot use atomic for max(in C/C++)
- May be better if starting/ending a parallel section is expensive

イロト イヨト イヨト イヨ

```
#pragma omp parallel
} ob
  #pragma omp single
    err = 0:
    k = 1 - k;
  }
  #pragma omp for collapse(2) reduction(max:err)
  for (i=1; i<M-1; i++)</pre>
    for (j=1; j<N-1; j++) {
      a[1-k][i][j] = 0.25 * (a[k][i-1][j] + a[k][i+1][j] +
                              a[k][i][j-1] + a[k][i][j+1]);
      err = fmax(err, fabs(a[1][i][j]-a[0][i][j]));
    }
  /* implicit barrier at end of omp for */
} while (err > maxerr);
```

< 日 > < 同 > < 三 > < 三 >

Metrics

Marc Snir

▶ ◀ 볼 ▶ 볼 ∽ ९ (~ Fall 2018 38 / 47

▲□▶ ▲圖▶ ▲厘▶ ▲厘▶

- Sequential performance $T_1(n)$: time to solve a problem of size n. (Same as W(n), computation work needed to solve the problem.)
 - To simplify, assume (wrongly) each operation takes one time unit.
- Parallel performance $T_p(n)$: time to solve a problem of size *n* with *p* hardware threads
- Ideal world: Can run p time faster than with one hardware thread; $T_p(n) = T_1(n)/p$. World is not ideal.

- Sequential performance $T_1(n)$: time to solve a problem of size n. (Same as W(n), computation work needed to solve the problem.)
 - To simplify, assume (wrongly) each operation takes one time unit.
- Parallel performance $T_p(n)$: time to solve a problem of size *n* with *p* hardware threads
- Ideal world: Can run p time faster than with one hardware thread; $T_p(n) = T_1(n)/p$. World is not ideal.
 - Parallel code does more work than sequential code (e.g., spawning threads)

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

- Sequential performance $T_1(n)$: time to solve a problem of size n. (Same as W(n), computation work needed to solve the problem.)
 - To simplify, assume (wrongly) each operation takes one time unit.
- Parallel performance $T_p(n)$: time to solve a problem of size n with p hardware threads
- Ideal world: Can run p time faster than with one hardware thread; $T_p(n) = T_1(n)/p$. World is not ideal.
 - Parallel code does more work than sequential code (e.g., spawning threads)
 - Parallel code may not have enough parallelism enough operations that can executed independently on distinct threads

▲□▶ ▲圖▶ ▲圖▶ ▲圖▶ ▲ 圖 - ∽○< ⊙

Marc Snir

・ロト ・ 日 ト ・ 日 ト ・ 日 ト

•
$$S_1(n) = 1$$

Normally S_p(n) ≤ p; there may be cases of superlinear speedup: p caches have larger capacity than one cache.

イロト イポト イヨト イヨト 一日

- $S_1(n) = 1$
- Normally S_p(n) ≤ p; there may be cases of superlinear speedup: p caches have larger capacity than one cache.
- Normally S_p(n) ≥ 1; but it is not rare to find that a parallel program runs more slowly than a sequential program solving the same problem (parallel overheads and lack of parallelism)

イロト イ団ト イヨト イヨト 二日

- $S_1(n) = 1$
- Normally S_p(n) ≤ p; there may be cases of superlinear speedup: p caches have larger capacity than one cache.
- Normally S_p(n) ≥ 1; but it is not rare to find that a parallel program runs more slowly than a sequential program solving the same problem (parallel overheads and lack of parallelism)

Efficiency Ratio between speedup and number of hardware threads: $E_p(n) = S_p(n)/p = T_1(n)/(p \times T_p(n))$; ratio between sequential work and parallel work. Normally $E_p(n) < 1$.

▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ 臣 のへの

```
#pragma omp parallel
    int i, j;
    time = omp_get_wtime();
    #pragma omp for collapse(2)
    for (i=0; i<N; i++)</pre>
      for (j=0; j<N; j++)</pre>
        c[i] += a[i][j] * b[j];
    printf("%d ",(int)(1000000.0 *
            (omp_get_wtime()-time)));
  }
  printf("\n");
}
```

"Fake" experiment: measuring running time of each thread and taking maximum

}

▲□▶ ▲圖▶ ▲ 圖▶ ▲ 圖▶ ― 圖 … のへで

 Execution time decreases as number of threads increase – But cannot decrease below the time to execute the inner loop (without code changes)

(I) < (II) < (II) < (II) < (II) < (II) < (III) </p>

• Superlinear speedup in some regions

Fall 2018 42 / 47

Theorem

If a computation has a fraction α that can be executed in parallel and a fraction $1-\alpha$ that is sequential, then

$$S_p = \frac{1}{(1-\alpha) + \alpha/p}$$

Proof.

$$T_p = \frac{\alpha T_1}{p} + (1 - \alpha) T_1$$
$$S_p = \frac{T_1}{T_p} = \frac{T_1}{\frac{\alpha T_1}{p} + (1 - \alpha) T_1} = \frac{1}{(1 - \alpha) + \alpha/p}$$

イロト イヨト イヨト イ

- Theorem holds as long as p does not exceed the level of available parallelism (N, in our example).
- Speedup can never exceed $\frac{1}{1-\alpha}$, the ratio between total work and sequential work.

▲□▶ ▲□▶ ▲三▶ ▲三▶ ▲□ ● ● ●

- T_1 number of operations
- T_{∞} longest critical path in the execution
- In our example ${\mathcal T}_1(N) \sim N^2$ and ${\mathcal T}_\infty(N) \sim N$
- Claim1: $T_p \ge \max(\frac{T_1}{p}, T_\infty)$
- Claim2: One can load-balance execution so that $T_p = O(\frac{T_1}{p} + T_\infty)$

・ロト ・聞 ト ・ 臣 ト ・ 臣 ト … 臣

Amdahl's Law revisited

- Amdahl's law was seen as evidence that parallelism has limited hope eventually, the sequential part dominates.
- What was missing is a realization that larger machines are used to solve larger problems
- *Strong scaling*: Consider a fixed size problem and apply more and more processors to its solution.
 - Eventually, more processors do not help. (Too many cooks spoil the broth.)
- Weak scaling: Increase problem size as the number of processors is increased – keep work per processor constant.
 - Ideally, total compute time stays fixed

- Strong scaling of Matrix×vector
- Usually scaling is sublinear because of overheads of communication and synchronization

イロト イヨト イヨト イヨト

How do we study this?

- Study a function in two variables: T_p(n) or S_p(n)
- Fix n and plot execution time T_p(n) as function of n. This is strong scaling: Keep the work fixed and increase the number of workers; efficiency decreases
- Keep work per processor fixed and study increase in compute time. This is *weak scaling*
- Keep Efficiency fixed and study how problem size increases:

 $n_{1/2}(p) = \min\{n: E_p(n) \ge 0.5\}$

Speedup as function of n and p

(日)

Fall 2018 47 / 47