Software Performance Engineering

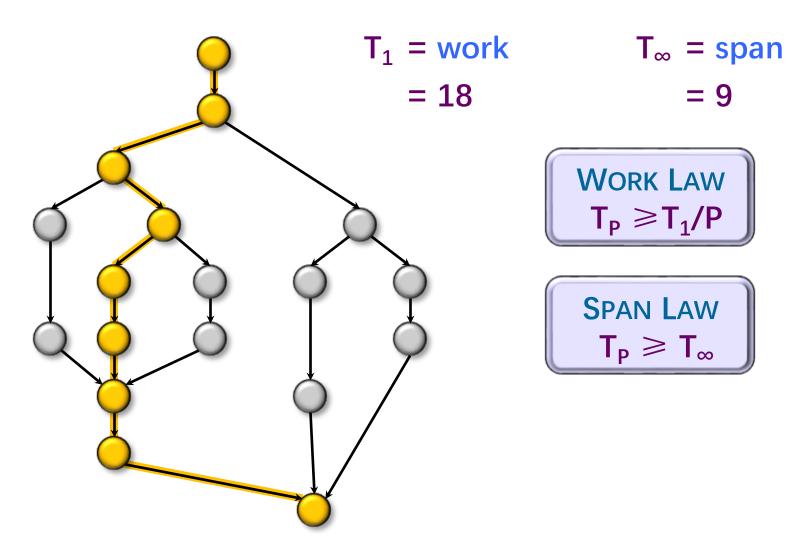
Scheduling Theory and Parallel Loops

Xuhao Chen October 14, 2025



Performance Measures

 T_P = execution time on P processors



Parallelism

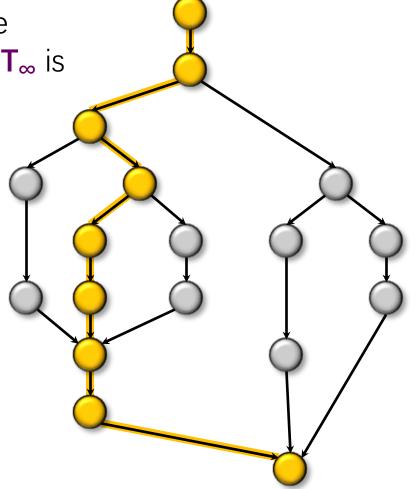
As the SPAN Law dictates that $T_P \ge T_{\infty}$, the maximum possible speedup given T_1 and T_{∞} is

 T_1/T_{∞} = parallelism

= the average amount of work per step along the span

= 18/9

= 2.

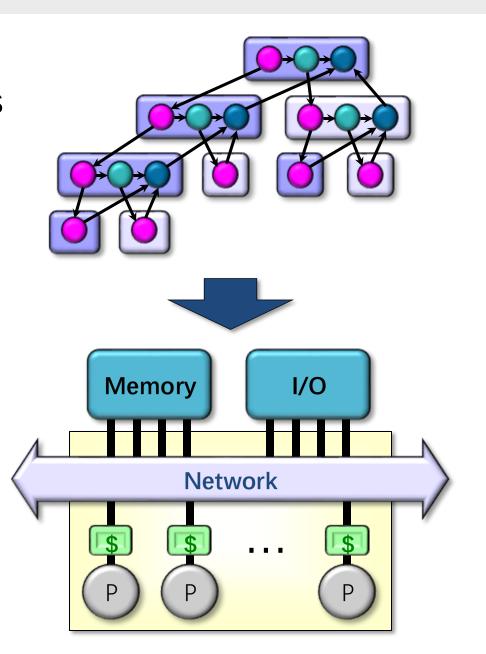




SCHEDULING THEORY

Scheduling

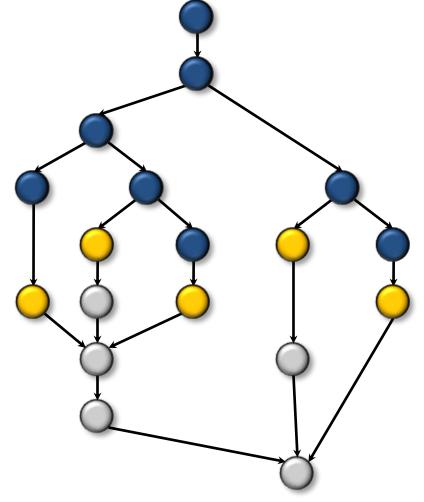
- Cilk allows the programmer to express potential parallelism in an application
- The Cilk scheduler maps strands onto processors dynamically at runtime
- Since the theory of distributed schedulers is complicated, we'll explore the ideas with a simple, centralized scheduler



Greedy Scheduling

IDEA: Do as much as possible on every step.

Definition. A strand is ready if all its predecessors have executed.



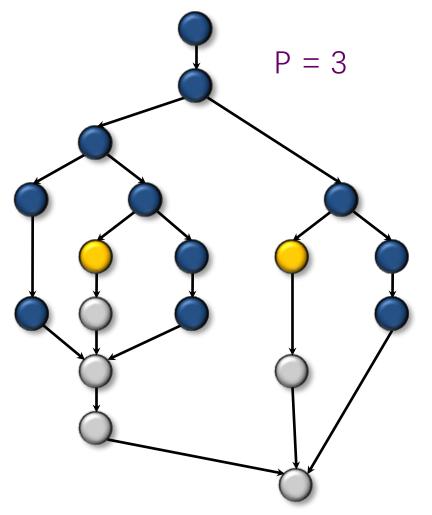
Greedy Scheduling

IDEA: Do as much as possible on every step.

Definition. A strand is ready if all its predecessors have executed.

Complete step

- ▶ P strands ready.
- Run any P.



Greedy Scheduling

IDEA: Do as much as possible on every step.

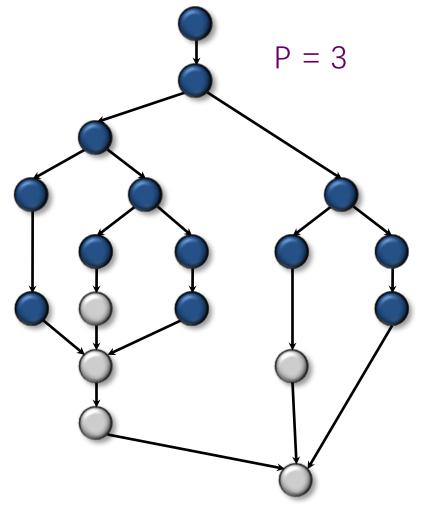
Definition. A strand is ready if all its predecessors have executed.

Complete step

- ▶ P strands ready.
- Run any P.

Incomplete step

- < P strands ready.
- Run all of them.



Analysis of Greedy

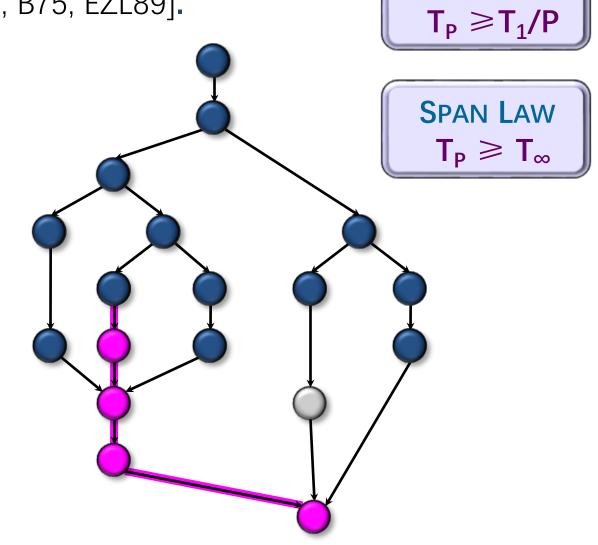
Greedy Scheduling Theorem [G68, B75, EZL89].

Any greedy scheduler achieves

$$T_{P} \leq T_{1}/P + T_{\infty}$$
.

Proof.

- # complete steps $\leq T_1/P$ since each complete step performs P work.
- # incomplete steps ≤ T_∞ since each incomplete step reduces the span of the unexecuted dag by 1.



WORK LAW

Optimality of Greedy

Corollary. Any greedy scheduler achieves within a factor of 2 of optimal.

Proof. Let T_P^* be the execution time produced by the optimal scheduler. Since $T_P^* \ge \max\{T_1/P, T_\infty\}$ by the WORK and SPAN LAWS, we have

$$T_P \le T_1/P + T_{\infty}$$
 $\le 2 \cdot \max\{T_1/P, T_{\infty}\}$
 $\le 2T_P *$

Linear Speedup

Corollary. Any greedy scheduler achieves near-perfect linear speedup whenever $T_1/T_\infty \gg P$.

Proof. Since $T_1/T_\infty \gg P$ is equivalent to $T_\infty \ll T_1/P$, the Greedy Scheduling Theorem gives us

$$T_P \leq T_1/P + T_{\infty}$$

 $\approx T_1/P$.

Thus, the speedup is $T_1/T_P \approx P$.

Definition. The quantity $(T_1/T_{\infty})/P = T_1/PT_{\infty}$ is called the parallel slackness.

Cilk Performance

- Cilk's randomized work-stealing scheduler achieves
 - $T_P = T_1/P + O(T_{\infty})$ expected time (provably);
 - \star T_P \approx T₁/P + T_{∞} time (empirically).
- Near-perfect linear speedup as long as $P \ll T_1/T_{\infty}$.
- Instrumentation in Cilkscale allows you to measure T_1 and T_{∞} .



PARALLEL LOOPS

Loop Parallelism in Cilk

Example:

In-place matrix transpose

The iterations of a cilk_for loop execute in parallel.

```
// indices run from 0, not 1
for (int i=1; i<n; ++i) {
  for (int j=0; j<i; ++j) {
    double temp = A[i][j];
    A[i][j] = A[j][i];
    A[j][i] = temp;
}
}</pre>
```

Loop Parallelism in Cilk

Example:

In-place matrix transpose

The iterations of a cilk_for loop execute in parallel.

```
// indices run from 0, not 1
cilk_for (int i=1; i<n; ++i) {
   for (int j=0; j<i; ++j) {
      double temp = A[i][j];
      A[i][j] = A[j][i];
      A[j][i] = temp;
   }
}</pre>
```

```
// indices run from 0, not 1
cilk_for (int i=1; i<n; ++i) {</pre>
  for (int j=0; j<i; ++j) {</pre>
    double temp = A[i][j]; void p_loop(int lo, int hi) //half open
    A[i][j] = A[j][i];
    A[j][i] = temp;
```

Original code

Divide-and-conquer

The OpenCilk compiler implements cilk_for loops using divide and conquer at optimization levels **-01** and higher

Compiler-generated recursion

```
if (hi > lo + 1) {
   int mid = lo + (hi - lo)/2;
    cilk scope {
      cilk_spawn p_loop(lo, mid);
      p_loop(mid, hi);
    return;
  int i = lo;
  for (int j=0; j<i; ++j) {
    double temp = A[i][j];
   A[i][j] = A[j][i];
   A[j][i] = temp;
p_loop(1, n);
```

Original code

Divide-and-conquer

The OpenCilk compiler implements cilk_for loops using divide and conquer at optimization levels -01 and higher

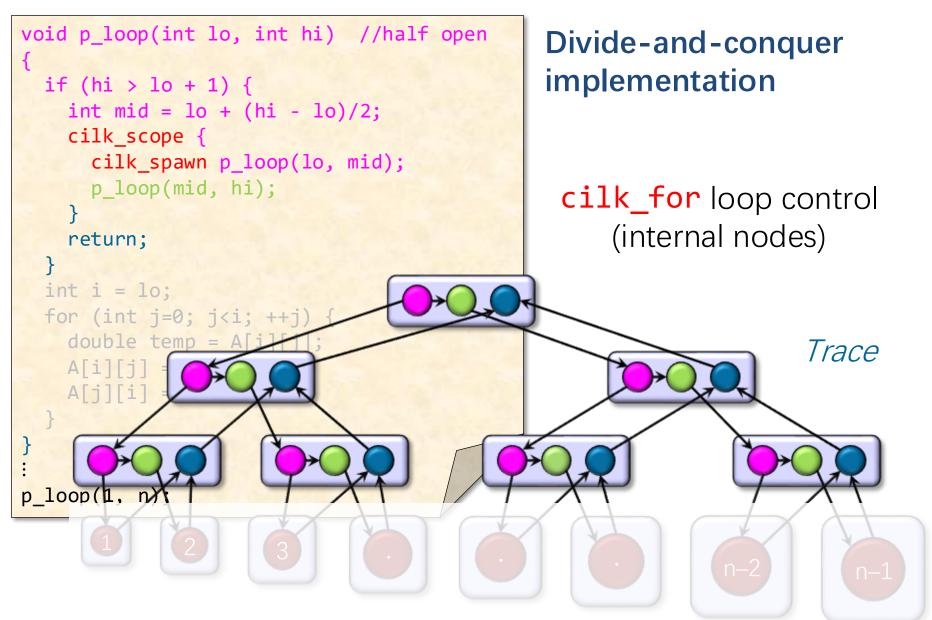
Compiler-generated recursion

```
void p_loop(int lo, int hi) //half open
  if (hi > lo + 1) {
    int mid = lo + (hi - lo)/2;
    cilk scope {
      cilk_spawn p_loop(lo, mid);
      p loop(mid, hi);
    return;
  int i = lo;
  for (int j=0; j<i; ++j) {
    double temp = A[i][j];
   A[i][j] = A[j][i];
    A[j][i] = temp;
p_loop(1, n);
```

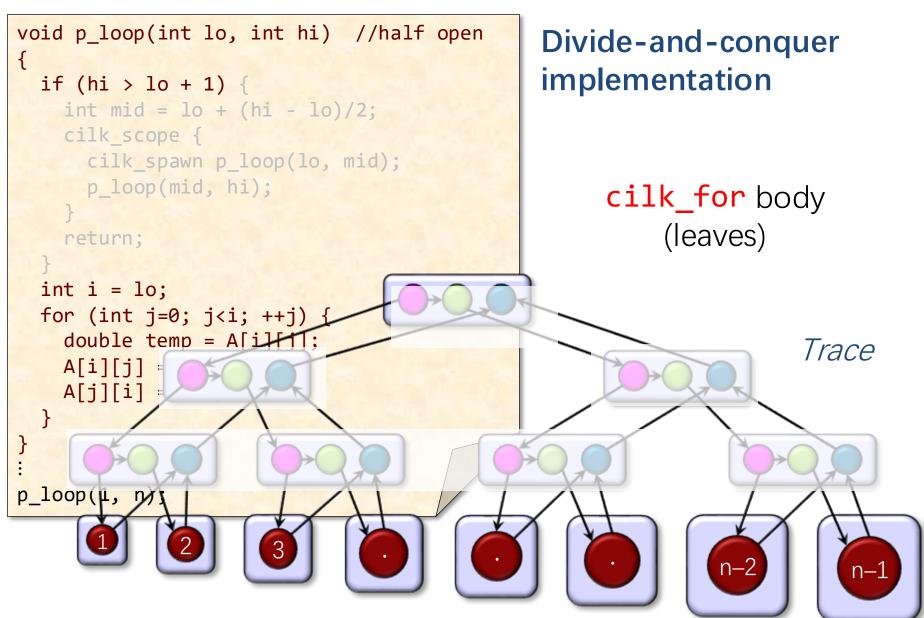
```
// indices run from 0, not 1
                                                        cilk for
cilk_for (int i=1; i<n; ++i) {</pre>
                                                       loop control
  for (int j=0; j<i; ++j) {
    double temp = A[i][j];[
                             void p_loop(int lo, int hi) //half open
    A[i][j] = A[j][i];
                               if (hi > lo + 1) {
    A[j][i] = temp;
                                 int mid = lo + (hi - lo)/2;
                                 cilk scope{
                                   cilk spawn p loop(lo, mid);
                                   p loop(mid, hi);
                                 return;
                               int i = lo;
                               for (int j=0; j<i; ++j) {
                                 double temp = A[i][j];
                                 A[i][j] = A[j][i];
                                 A[j][i] = temp;
                             p_loop(1, n);
```

```
// indices run from 0, not 1
cilk_for (int i=1; i<n; ++i) {</pre>
  for (int j=0; j<i; ++j) {</pre>
    double temp = A[i][j];
                              void p_loop(int lo, int hi) //half open
    A[i][j] = A[j][i];
    A[j][i] = temp;
                              if (hi > lo + 1) {
                                  int mid = lo + (hi - lo)/2;
                                  cilk scope {
                                    cilk_spawn p_loop(lo, mid);
                                    p_loop(mid, hi);
                                  return;
                                int i = lo;
                                for (int j=0; j<i; ++j) {</pre>
                                  double temp = A[i][j];
                                  A[i][j] = A[j][i];
                                  A[j][i] = temp;
           lifted
         loop body
                              p_loop(1, n);
```

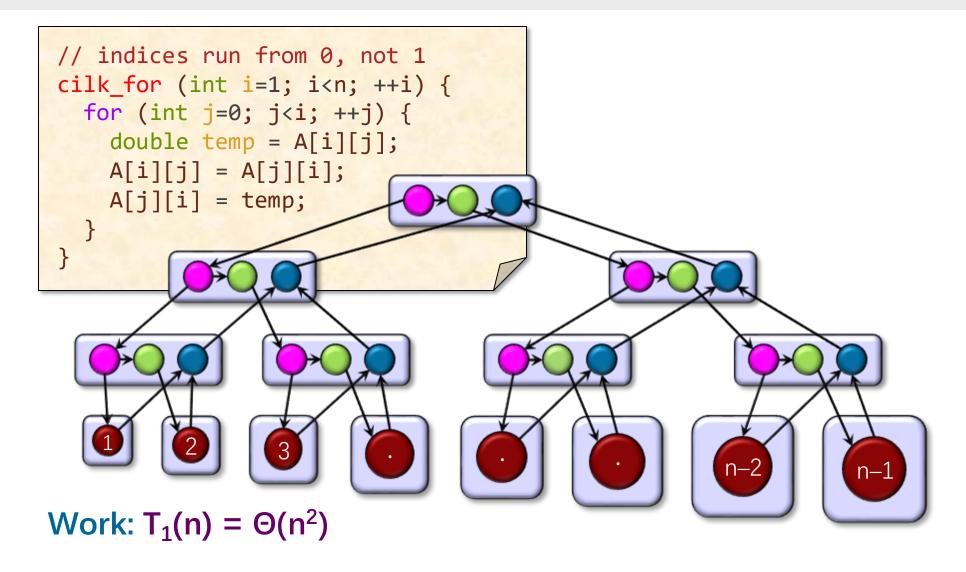
Execution of Parallel Loops



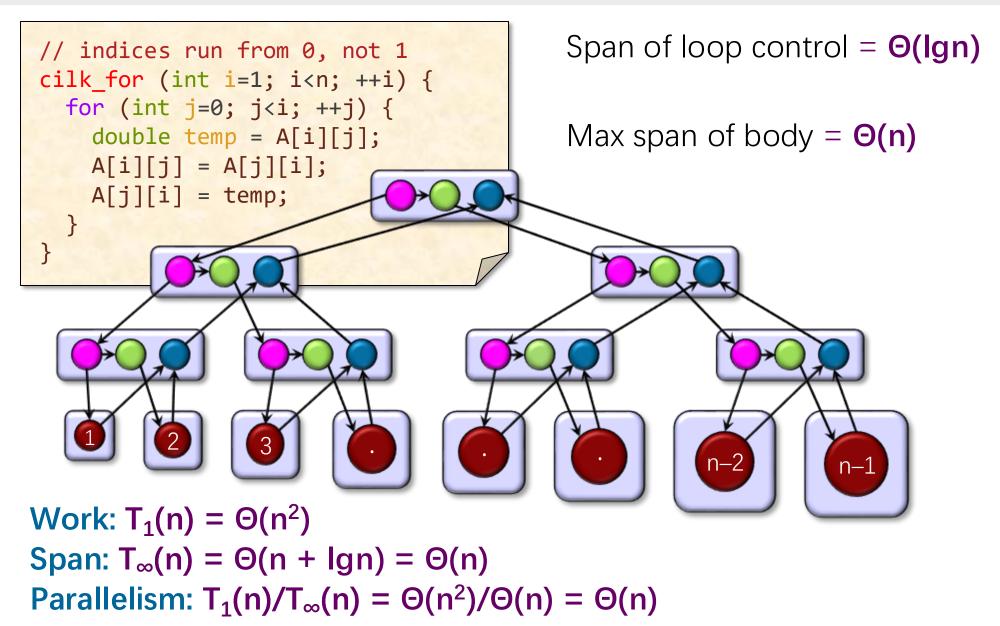
Execution of Parallel Loops



Analysis of Parallel Matrix Transpose



Analysis of Parallel Matrix Transpose



Analysis of Nested Parallel Loops

```
// indices run from 0, not 1
cilk_for (int i=1; i<n; ++i) {
   cilk_for (int j=0; j<i; ++j) {
      double temp = A[i][j];
      A[i][j] = A[j][i];
      A[j][i] = temp;
   }
}</pre>
```

Span of outer loop control = $\Theta(lgn)$

Max span of inner loop control = $\Theta(lgn)$

Span of body = $\Theta(1)$.

```
Work: T_1(n) = \Theta(n^2)
Span: T_{\infty}(n) = \Theta(\lg n)
Parallelism: T_1(n)/T_{\infty}(n) = \Theta(n^2/\lg n)
```

Analysis of Nested Parallel Loops

```
// indices run from 0, not 1
cilk_for (int i=1; i<n; ++i) {
   cilk_for (int j=0; j<i; ++j) {
      double temp = A[i][j];
      A[i][j] = A[j][i],
      A[j][i] = temp;
   }
}</pre>
```

Span of outer loop control = $\Theta(lgn)$

Max span of inner loop control = $\Theta(lgn)$

of body = $\Theta(1)$.

How much loop control overhead?

Work: $T_1(n) = \Theta(n^2)$

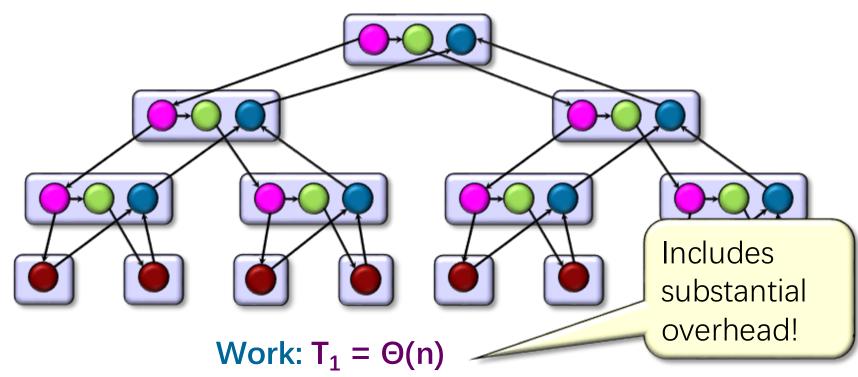
Span: $T_{\infty}(n) = \Theta(Ign)$

Parallelism: $T_1(n)/T_{\infty}(n) = \Theta(n^2/\lg n)$

A Closer Look at Parallel Loops

Vector addition

```
cilk_for (int i=0; i<n; ++i) {
   A[i] += B[i];
}</pre>
```



Span: $T_{\infty} = \Theta(Ign)$

Parallelism: $T_1/T_{\infty} = \Theta(n/\lg n)$

Optimizing Parallel-Loop Control

```
cilk_for (int i=0; i<n; ++i) {
   A[i] += B[i];
}</pre>
```

Original code

Compiler-generated recursion

```
void p_loop(int lo, int hi) { //half open
  if (hi > lo + 1) {
    int mid = lo + (hi - lo)/2;
    cilk_scope {
      cilk_spawn p_loop(lo, mid);
      p_loop(mid, hi);
    return;
  for (int i=lo; i<hi; ++i) {</pre>
    A[i] += B[i];
p_loop(0, n);
```

Coarsening Parallel Loops

```
#pragma cilk grainsize G
cilk_for (int i=0; i<n; ++i) {
   A[i] += B[i];
}</pre>
```

If a grain-size pragma is not specified, the Cilk runtime system heuristically guesses

G to minimize overhead.

Compiler-generated recursion

```
void p_loop(int lo, int hi) { //half open
  if (hi > lo + G) {
    int mid = lo + (hi - lo)/2;
    cilk scope {
      cilk spawn p loop(lo, mid);
      p_loop(mid, hi);
    return;
  for (int i=lo; i<hi; ++i) {</pre>
    A[i] += B[i];
p_loop(0, n);
```

#pragma cilk grainsize G Vector cilk_for (int i=0; i<n; ++i) {</pre> addition A[i] += B[i]; $G \cdot I$

Let I be the time for one iteration of the loop body.

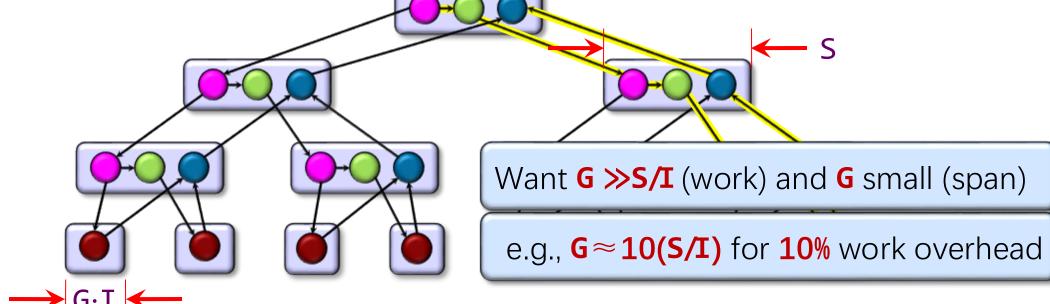
Let **S** be the time to perform a level of the recursion.

#pragma cilk grainsize G Vector cilk_for (int i=0; i<n; ++i) {</pre> addition A[i] += B[i];Work: $T_1 = n \cdot I + (n/G - 1) \cdot S$ real work work overhead

#pragma cilk grainsize G Vector cilk_for (int i=0; i<n; ++i) {</pre> addition A[i] += B[i];G be large Work: $T_1 = n \cdot I + (n/G - 1) \cdot S$ Span: $T_{\infty} = G \cdot I + \frac{1}{16} \frac{(n/G) \cdot S}{1}$ G be small

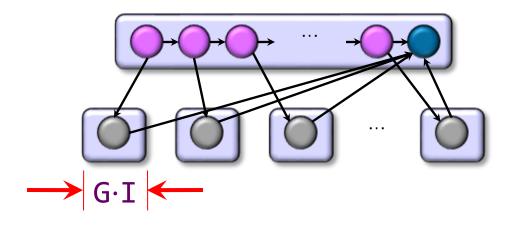
#pragma cilk grainsize G Vector cilk_for (int i=0; i<n; ++i) {</pre> addition A[i] += B[i];Work: $T_1 = n \cdot I + (n/G - 1) \cdot S$ **Span:** $T_{\infty} = G \cdot I + \lg(n/G) \cdot S$

#pragma cilk grainsize G
cilk_for (int i=0; i<n; ++i) {
 A[i] += B[i];
}</pre>



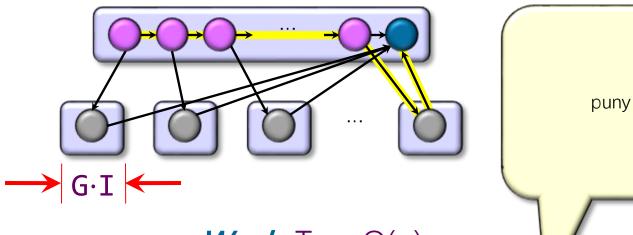
Work: $T_1 = n \cdot I + (n/G - 1) \cdot S$ Span: $T_{\infty} = G \cdot I + \lg(n/G) \cdot S$ Parallelism $T_1/T_{\infty} \approx \Theta(n/lgn)/(S/I)$

Another Implementation



Work: $T_1 = \Theta(n)$ Assume that G = 1. Span: $T_{\infty} = 0$

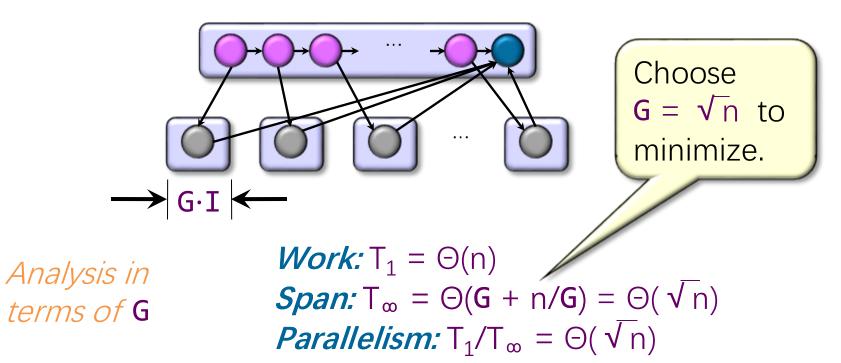
Another Implementation



Assume that G = 1.

Work: $T_1 = \Theta(n)$ Span: $T_{\infty} = \Theta(n)$ Parallelism: $T_1/T_{\infty} = \Theta(1)$

Another Implementation



Quiz on Parallel Loops

Question: Let $P \ll n$ be the number of workers on the system. How does the asymptotic parallelism of Code A compare to that of Code B? (Differences highlighted.)

Code A

```
#pragma cilk grainsize 1
cilk_for (int i=0; i<n; i+=32) {
  for (int j=i; j<MIN(i+32, n); ++j)
        A[j] += B[j];
}</pre>
```

Work:

Span:

Parallelism:

Code B

```
#pragma cilk grainsize 1
cilk_for (int i=0; i<n; i+=n/P) {
  for (int j=i; j<MIN(i+n/P, n); ++j)
        A[j] += B[j];
}</pre>
```

Work:

Span:

Parallelism:

Three Performance Tips

- 1. Minimize the span to maximize parallelism.
 - Try to generate **10** times more parallelism than processors for near-perfect linear speedup.
- 2. If you have plenty of parallelism, try to trade some of it off to reduce work overhead.
- 3. Use divide-and-conquer recursion or parallel loops rather than spawning one small thing after another.

And Three More

- 4. Ensure that work/#spawns is sufficiently large.
 - Coarsen by using function calls and inlining near the leaves of recursion, rather than spawning.
- 5. Parallelize **outer loops**, as opposed to inner loops, if you're forced to make a choice.
- 6. Watch out for scheduling overheads.

```
Do this:
```

```
cilk_for (int i=0; i<2; ++i) {
  for (int j=0; j<n; ++j)
   f(i,j);
}</pre>
```

Not this:

```
for (int j=0; j<n; ++j) {
   cilk_for (int i=0; i<2; ++i)
     f(i,j);
}</pre>
```



Take-Aways



- Any greedy scheduler provides linear speedup on computations having sufficient parallel slackness
- The OpenCilk runtime system incorporates a randomized work-stealing scheduler that has strong theoretical bounds on its running time which are similar to those for greedy scheduling
- Loops in Cilk are synthesized using divide-and-conquer spawning, which incurs linear work and logarithmic span
- Coarsening recursion can reduce loop overhead