

# Parallelization

Rupesh Nasre.

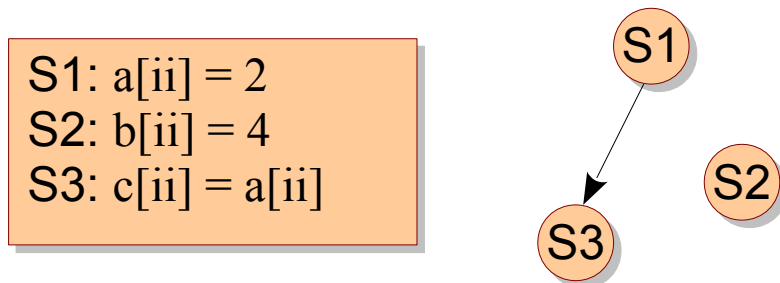
CS6843 Program Analysis  
IIT Madras  
Jan 2016

# Speedup

- Speedup =  $T_s / T_p$
- **Amdahl's Law:** Speedup is limited by the sequential part of the task.
- If 20% of the task is sequential, program's speedup is limited to 5 (irrespective of the number of cores or amount of effort).

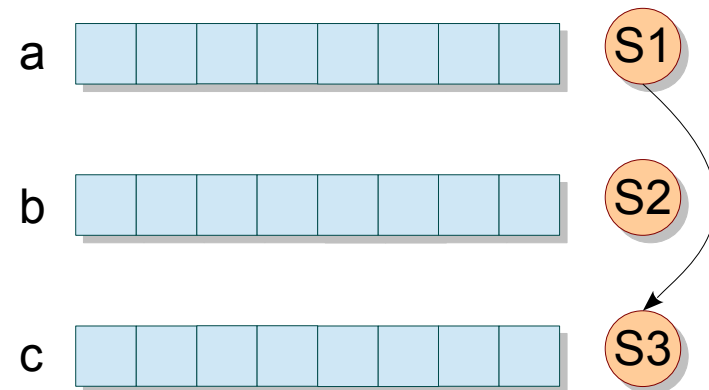
# Instruction Parallel vs. Data Parallel

- Parallelism extracted from multiple instructions on the data items.



S1 is the source and S3 is the sink of the dependence.

- Parallelism extracted from the same task on different data items.



# Control Dependence

- `if (x == 4) y = 10; else y = 1;`

# Data Dependence

- $\text{pi} = 3.142; r = 5.0; \text{area} = \text{pi} * r * r;$
- Types
  - True / Flow / RAW:  $S1 \delta S2 (x = \dots; \dots = x;)$
  - Anti / WAR:  $S1 \delta^{-1} S2 (\dots = x; x = \dots)$
  - Output / WAW:  $S1 \delta^0 S2 (x = \dots; x = \dots)$

# Program Order vs. Dependence

- **Sequential** order imposed by the program is too restrictive.
- Only the **partial order** of all dependences need to be maintained by the compiler to guarantee program correctness.
- So, reorder flow; maintain dependence.

# Advantages of Reordering

- Improved **locality**
  - Spatial: matrix operations
  - Temporal: `xinit(); yinit(); xcompute(); ycompute();`
- Improved **load balance**
  - `small1(); big1(); small2(); big2();`
- Improved **parallelism**
  - `xuse(); xdef(); yuse(); ydef();`

# Let's Focus on Loops

- **Iteration vector:** Sequence of outer loops.
  - $\vec{iv} = (i_{\text{outermost}}, \dots, i_{\text{middle}}, \dots, i_{\text{innermost}})$
  - For instance  $(i, j, k)$ .
- **Iteration space:** Set of all possible iteration vectors for a statement.
- **Statement instance:**  $S(\vec{i})$
- $S(\vec{i}) \delta S(\vec{j})$  iff
  - (a)  $i < j$  or  $(i == j \text{ and } S1 \Rightarrow \Rightarrow \Rightarrow S2 \text{ path in loop-body})$
  - (b) both access the same memory location
  - (c) at least one of the accesses is a write



# Safe Transformations

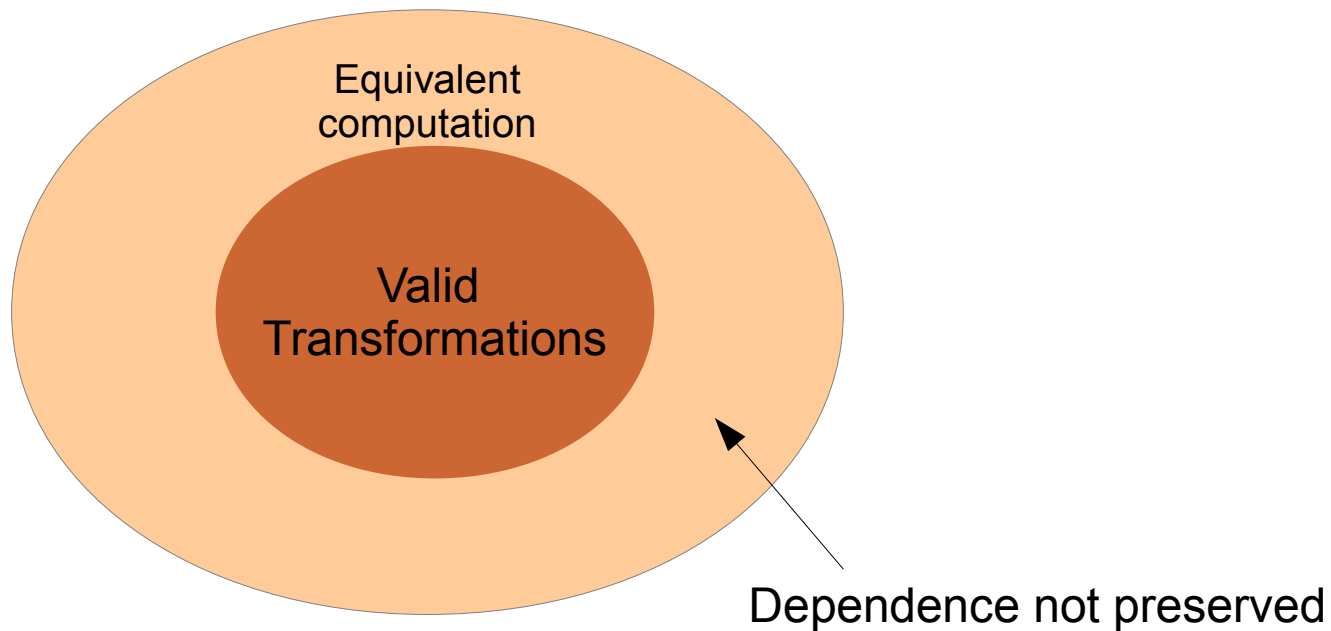
- **Loop Dependence Theorem**
  - There exists a dependence from statement  $S1$  to statement  $S2$  in a common nest of loops iff there exist two iteration vectors  $\vec{i}$  and  $\vec{j}$  for the nest, such that  $S1(\vec{i}) \delta S2(\vec{j})$ .
- Two computations are **equivalent** if on the same inputs they produce the same output.
- A transformation is **safe** if it leads to an equivalent program.

# Reordering Transformations

- A **reordering transformation** is any program transformation that merely changes the execution order of the code, without adding or deleting any executions of any statements.
- A reordering transformation preserves a dependence if it preserves the relative execution order of the source and the sink of that dependence.
- **Theorem:** Any reordering transformation that preserves every dependence in a program leads to an equivalent computation.

# Valid Transformations

- A transformation is valid for the program to which it applies if it preserves all the dependences in the program.



**Classwork:** Write a simple transformation that maintains computation equivalence but does not preserve dependence.

# Loop Parallelization

- **Theorem:** It is valid to convert a sequential loop to a parallel loop if the loop carries no dependence.

```
for (k = 0; k < n; ++k) {  
  S1: a[k] = b[k];  
  S2: b[k] = a[k] + 1;  
}
```



```
for (k = 0; k < n; ++k) {  
  S1: a[k] = a[k + 1];  
}
```



# General Strategy

```
for (ii = 0; ii < n; ++ii) {  
  for (jj = 0; jj < m; ++jj) {  
    a[f(ii, jj)][g(ii, jj)] = ...  
    ... = ... a[h(ii, jj)][k(ii, jj)]...  
  }  
}
```

Conditions for flow dependence from iteration  $(ii_w, jj_w)$  to  $(ii_r, jj_r)$ :

$$0 \leq ii_w < n$$

$$0 \leq jj_w < m$$

$$0 \leq ii_r < n$$

$$0 \leq jj_r < m$$

$$(ii_w, jj_w) \leq (ii_r, jj_r)$$

$$f(ii_w, jj_w) = h(ii_r, jj_r)$$

$$g(ii_w, jj_w) = k(ii_r, jj_r)$$

If  $f$ ,  $g$ ,  $h$ ,  $k$  are affine functions of loop variables, then dependence testing can be formulated as an ILP.

# ILP Formulation

```
for (ii = 0; ii < 10; ++ii) {  
    a[2 * ii] = ... a[2 * ii + 1] ...  
}
```

Is there a flow dependence between different iterations?

Dependence equations

$$0 \leq ii_w < ii_r < 10$$

$$2 * ii_w = 2 * ii_r + 1$$

which can be written as

$$0 \leq ii_w$$

$$ii_w \leq ii_r - 1$$

$$ii_r \leq 9$$

$$2 * ii_w \leq 2 * ii_r + 1$$

$$2 * ii_r + 1 \leq 2 * ii_w$$

$$\begin{pmatrix} -1 & 0 \\ 1 & -1 \\ 0 & 1 \\ 2 & -2 \\ -2 & 2 \end{pmatrix} \begin{pmatrix} ii_w \\ ii_r \end{pmatrix} \leq \begin{pmatrix} 0 \\ -1 \\ 9 \\ 1 \\ -1 \end{pmatrix}$$

Dependence exists if the system has a solution.

# ILP Formulation

```
for (ii = 0; ii < 10; ++ii) {  
    a[2 * ii] = ... a[2 * ii + 1] ...  
}
```

Is there an anti-dependence between different iterations?

Dependence equations

$$0 \leq ii_r < ii_w < 10$$

$$2 * ii_w = 2 * ii_r + 1$$

which can be written as

$$0 \leq ii_r$$

$$ii_r \leq ii_w - 1$$

$$ii_w \leq 9$$

$$2 * ii_w \leq 2 * ii_r + 1$$

$$2 * ii_r + 1 \leq 2 * ii_w$$

$$\begin{pmatrix} -1 & 0 \\ 1 & -1 \\ 0 & 1 \\ -2 & 2 \\ 2 & -2 \end{pmatrix} \begin{pmatrix} ii_r \\ ii_w \end{pmatrix} \leq \begin{pmatrix} 0 \\ -1 \\ 9 \\ 1 \\ -1 \end{pmatrix}$$

The system is not satisfiable, so anti-dependence does not exist.

# ILP Formulation

```
for (ii = 0; ii < 10; ++ii) {  
    a[2 * ii] = ... a[ii + 1] ...  
}
```

Is there an anti-dependence between different iterations?

Dependence equations

$$0 \leq ii_r < ii_w < 10$$

$$2 * ii_w = ii_r + 1$$

which can be written as

$$0 \leq ii_r$$

$$ii_r \leq ii_w - 1$$

$$ii_w \leq 9$$

$$2 * ii_w \leq ii_r + 1$$

$$ii_r + 1 \leq 2 * ii_w$$

$$\begin{pmatrix} -1 & 0 \\ 1 & -1 \\ 0 & 1 \\ -1 & 2 \\ 1 & -2 \end{pmatrix} \begin{pmatrix} ii_r \\ ii_w \end{pmatrix} \leq \begin{pmatrix} 0 \\ -1 \\ 9 \\ 1 \\ -1 \end{pmatrix}$$

The system is not satisfiable, so anti-dependence does not exist.



# ILP Formulation

```
for (ii = 0; ii < 10; ++ii) {
    a[2 * ii] = ... a[ii + 1] ...
}
```

Is there a true dependence between different iterations?

Dependence equations

$$0 \leq ii_w < ii_r < 10$$

$$2 * ii_w = ii_r + 1$$

which can be written as

$$0 \leq ii_w$$

$$ii_w \leq ii_r - 1$$

$$ii_r \leq 9$$

$$2 * ii_w \leq ii_r + 1$$

$$ii_r + 1 \leq 2 * ii_w$$

$$\begin{pmatrix} 0 & -1 \\ -1 & 1 \\ 1 & 0 \\ -1 & 2 \\ 1 & -2 \end{pmatrix} \begin{pmatrix} ii_r \\ ii_w \end{pmatrix} \leq \begin{pmatrix} 0 \\ -1 \\ 9 \\ 1 \\ -1 \end{pmatrix}$$

$ii_r$	$ii_w$
0	--
1	--
2	--
3	2
4	--
5	3
6	--
7	4
8	--
9	5

The system is satisfiable, so true dependence exists.

# ILP Formulation

```
for (ii = 0; ii < 10; ++ii) {  
  a[2 * ii] = ... a[ii + 1] ...  
  A[3 + ii] = ... a[5 * ii] ...  
}
```

Is there a true dependence between different iterations?

We will have to model equations across all inter-iteration pairs of reads/writes.

- $2 * ii$  and  $ii + 1$
- $3 + ii$  and  $5 * ii$
- $2 * ii$  and  $5 * ii$
- $3 + ii$  and  $ii + 1$

How about  $2 * ii$  and  $3 + ii$ ?

How about  $ii + 1$  and  $5 * ii$ ?

If any of the systems is satisfiable, then true dependence exists.

# Managing Races


- Data-race between iterations  $p$  and  $q$  for element  $a[f(i)]$ .
- Critical section
  - Locks
  - Atomics
  - Barriers

# Inserting Locks

- Data-race between iterations p and q for element  $a[f(i)]$ .

```
if (i == p || i == q) {  
    lock(f(i));  
    ... perform operation ...  
    unlock(f(i));  
}
```

This operation could be same or different for the involved threads.



- e.g., Producer-consumer

```
produce() {  
    while (...) {  
        items.add(...);  
    }  
}
```

```
consume() {  
    e = items.remove();  
}
```

# Inserting Locks

- For multiple data items  $a[f(i)]$  and  $a[g(i)]$ 
  - Single lock
  - Multiple locks
- Multiple locks may lead to deadlock
  - may allow deadlock if it improves parallelism
- Deadlock avoidance may lead to livelock
  - may allow livelock if rare

# Inserting Locks

- Sometimes, a lock may be for a simple operation

```
if (i == p || i == q) {  
    lock(f(i));  
    sum += a[i];  
    unlock(f(i));  
}
```

- A simple critical section may be convertible to atomics.

# Inserting Atomics

- If the operation is simple
  - Primitive type
  - Single element
  - Relative update / read-write
- Example
  - Producer-consumer with single element update
- Types
  - increment, decrement
  - add, sub
  - min, max
  - `exch`, `CAS`

# Inserting Atomics

- **Classwork:** convert the following example from locks to atomics

```
if (i == p || i == q) {  
    lock(f(i));  
    sum += a[i];  
    unlock(f(i));  
}
```

- **Classwork:** write parallel slist insertion and deletion routines using atomics
- **Homework:** write parallel dlist insertion routine using atomics



# Inserting Locks

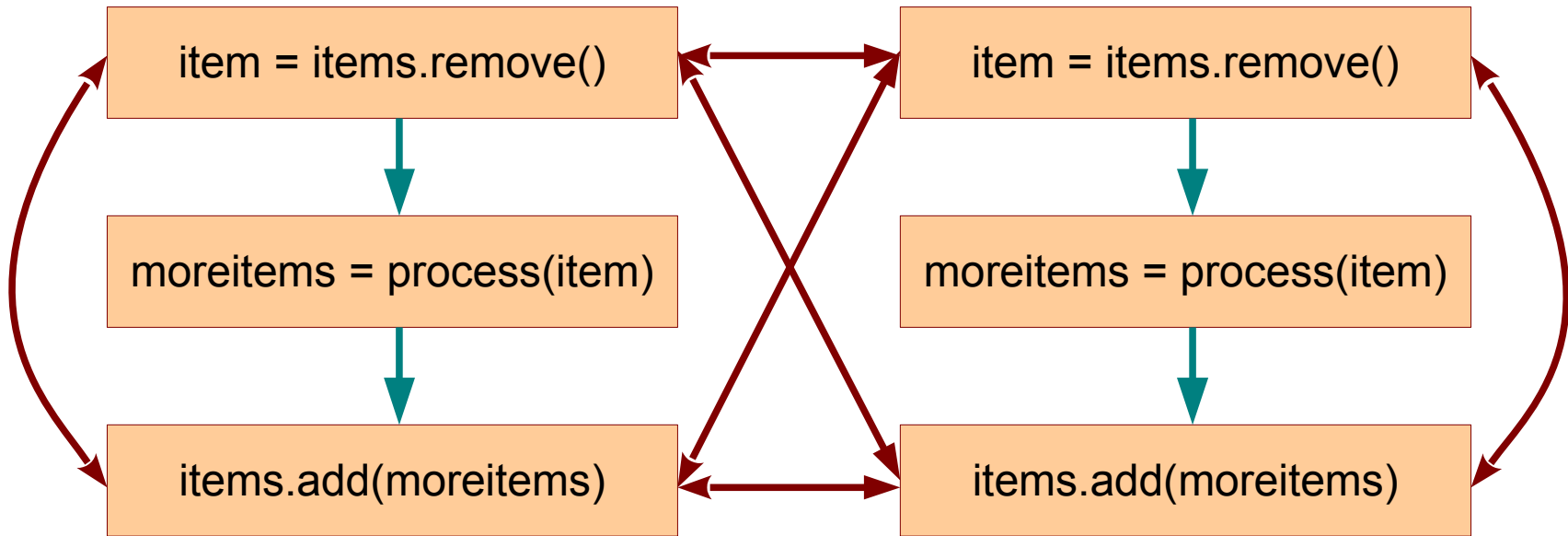
```
if (i == 1 || i == 2 || i == 4 || ...) {  
    lock(f(i));  
    item = items.remove();  
    moreitems = process(item);  
    items.add(moreitems);  
    unlock(f(i));  
}
```

- If there are many threads involved in the `if(...)` condition and the operation is multi-step, overapproximate the dependences.

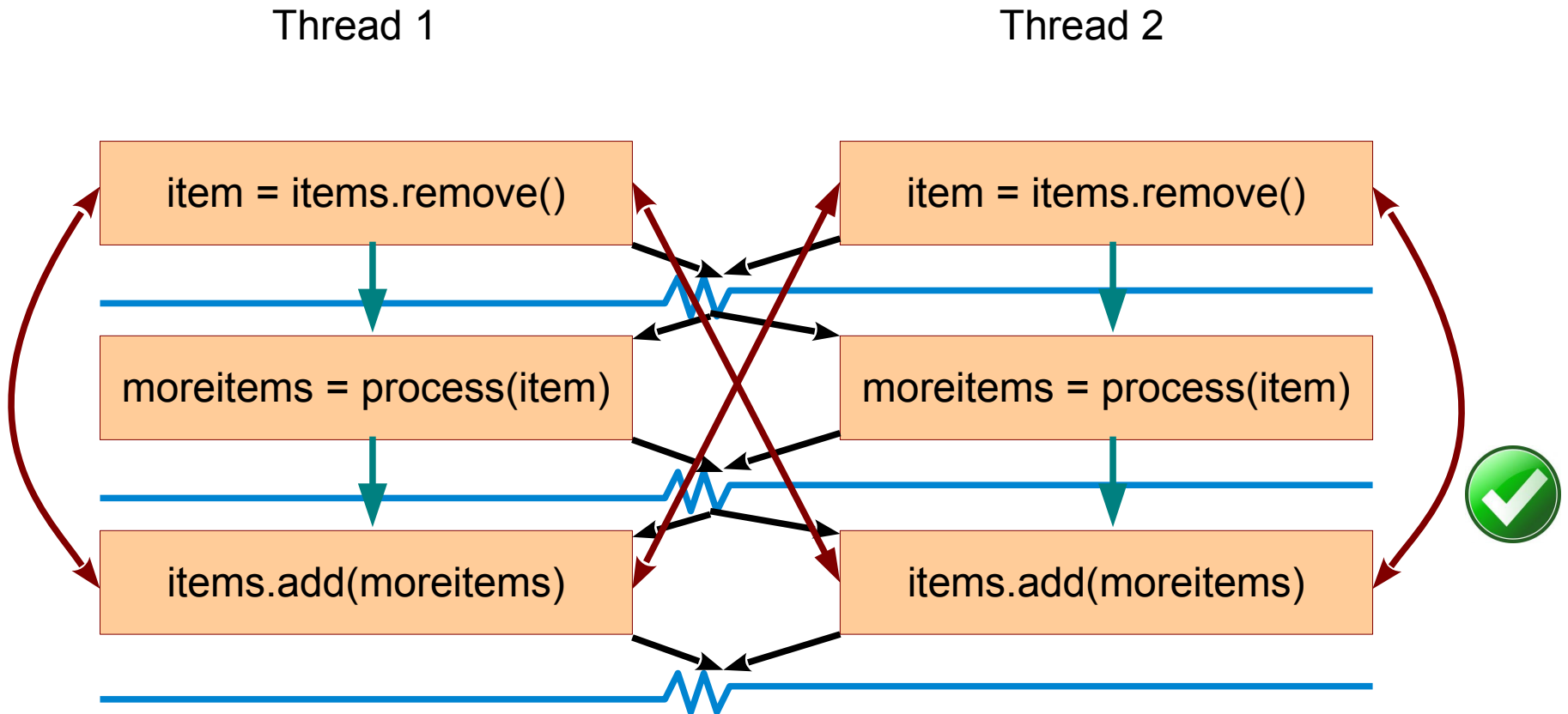
# Dependences

Thread 1

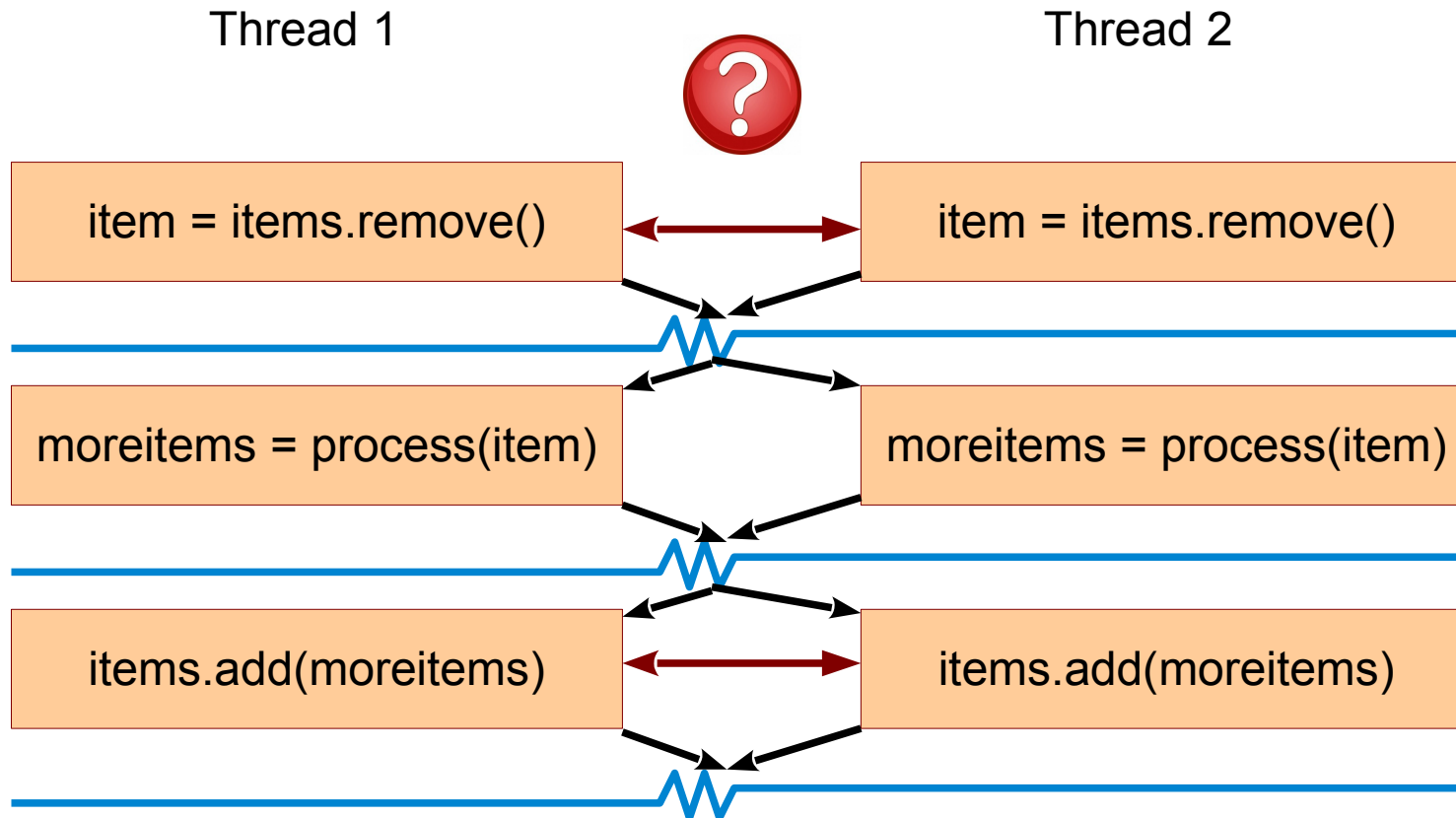
Thread 2



# Barriers



# Barriers



# Inserting Barriers

```
if (i == 1 || i == 2 || i == 4 || ...) {  
    lock(f(i));  
    item = items.remove();  
    unlock(f(i));  
    -- barrier --  
  
    moreitems = process(item);  
    -- barrier --  
  
    lock(f(i));  
    items.add(moreitems);  
    unlock(f(i));  
    -- barrier --  
}
```



Can be converted to atomics.



Can lead to good parallelism.



Can be converted to atomics.

# Inserting Barriers

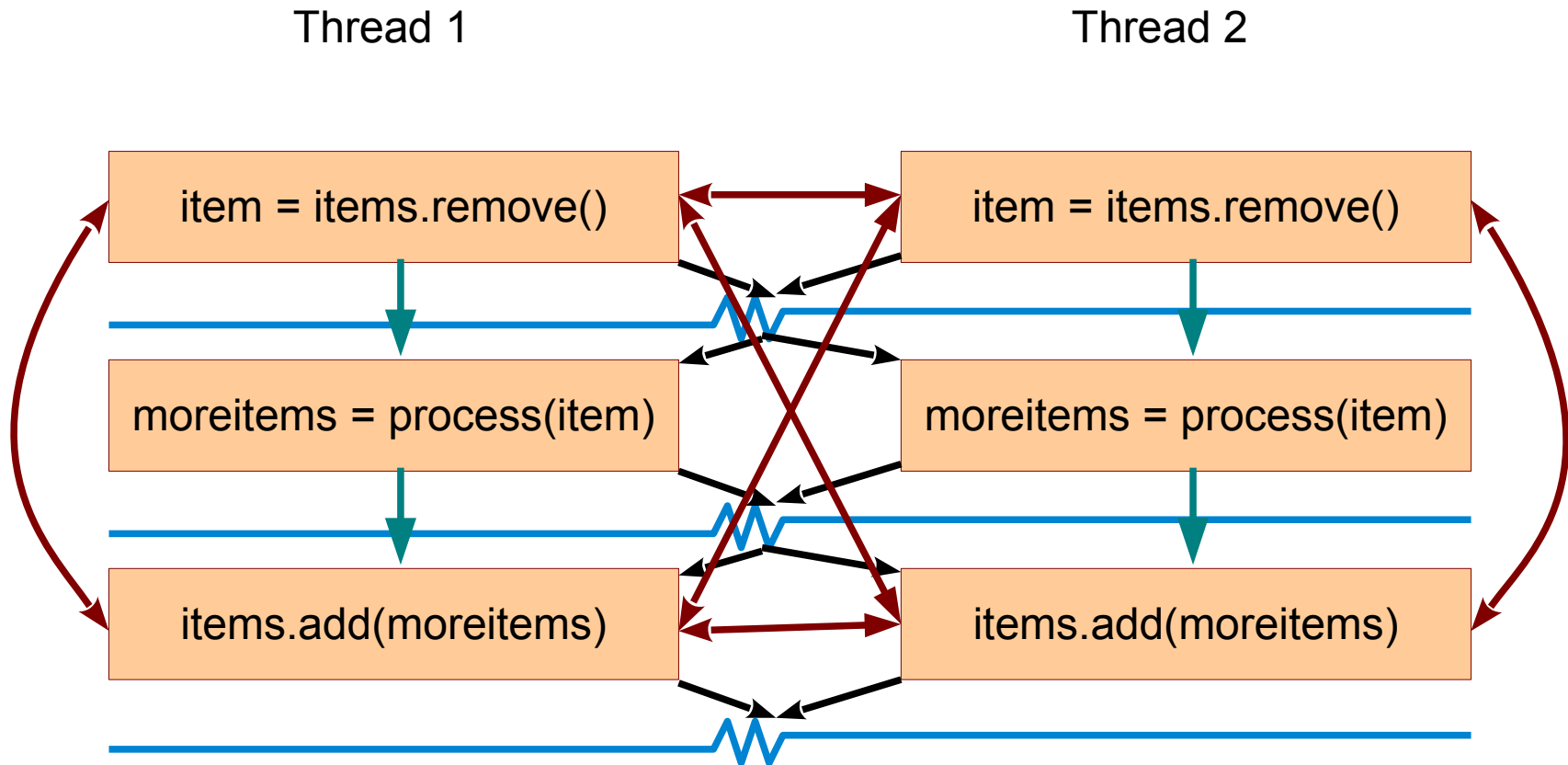
```
if (i == 1 || i == 2 || i == 4 || ...) {  
    atomicDec(items[f(i)]);  
    -- barrier --  
  
    moreitems = process(item);  
    -- barrier --  
  
    atomicAdd(items[f(i)], size(moreitems));  
    items.addunsync(moreitems);  
    -- barrier --  
}
```

If the barrier is emulated, one can combine these operations.

# Barriers and Dependences

- A barrier may be considered in effect similar to loop distribution.
- If dependences are sparse, use atomics/locks; otherwise barriers work well.
- A barrier may add more dependences than required.
- But it must preserve all the existing dependences.

# Barriers and Dependences



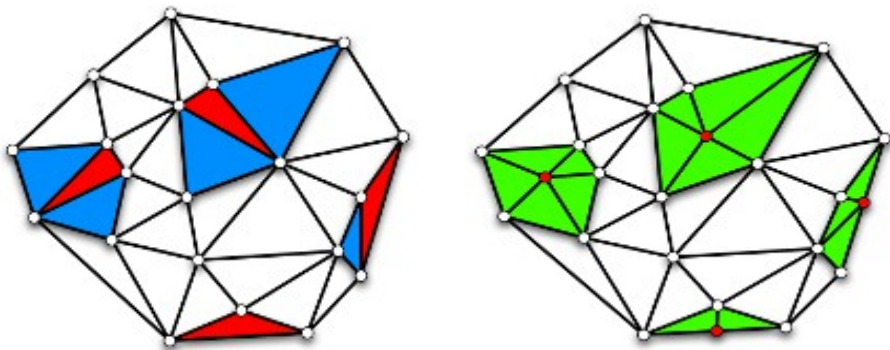
Did we add any extra dependences?



# Limitations of Static Parallelization

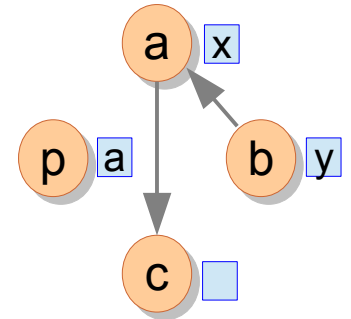
- Some programs cannot be effectively parallelized using static techniques.
  - e.g. graph algorithms, pointer-savvy programs
- Existing static optimization techniques (analysis) are also very conservative for such programs.
- Ineffectiveness of static techniques forces us to use dynamic approaches.

# Examples of Graph Algorithms

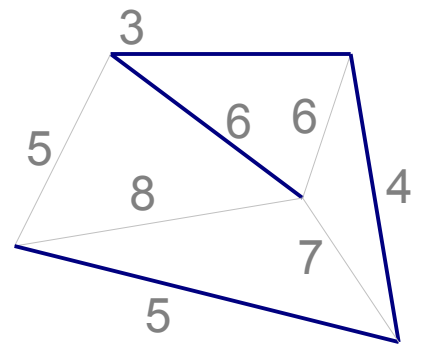


Delaunay Mesh Refinement

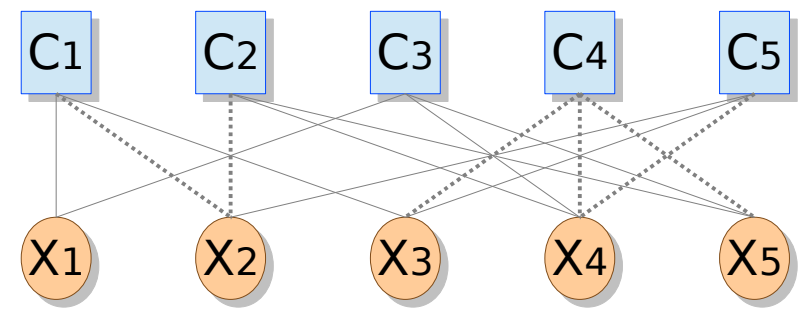
```
a = &x  
b = &y  
p = &a  
*p = b  
c = a
```



Points-to Analysis



Minimum Spanning Tree Computation



Survey Propagation

# What is IrRegularity?

- Data-access or control patterns are unpredictable at compile time.

## Irregular data-access

```
int a[N], b[N], c[N];  
readinput(a);  
  
c[5] = b[a[4]];
```

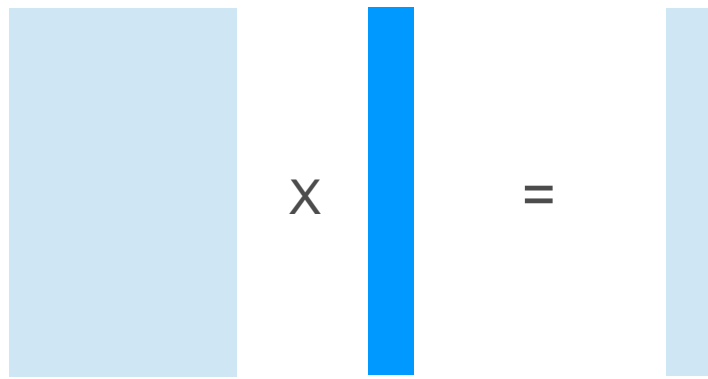
## Irregular control-flow

```
int a[N];  
readinput(a);  
  
if (a[4] > 30) {  
    ...  
}
```

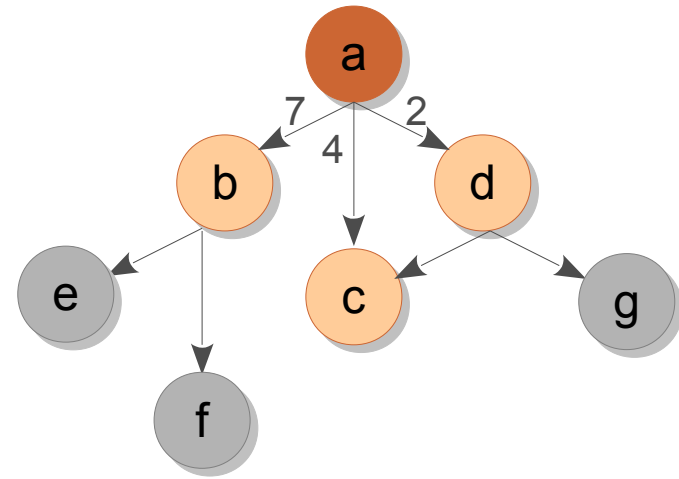
Needs dynamic techniques

Pointer-based data structures often contribute to irregularity.

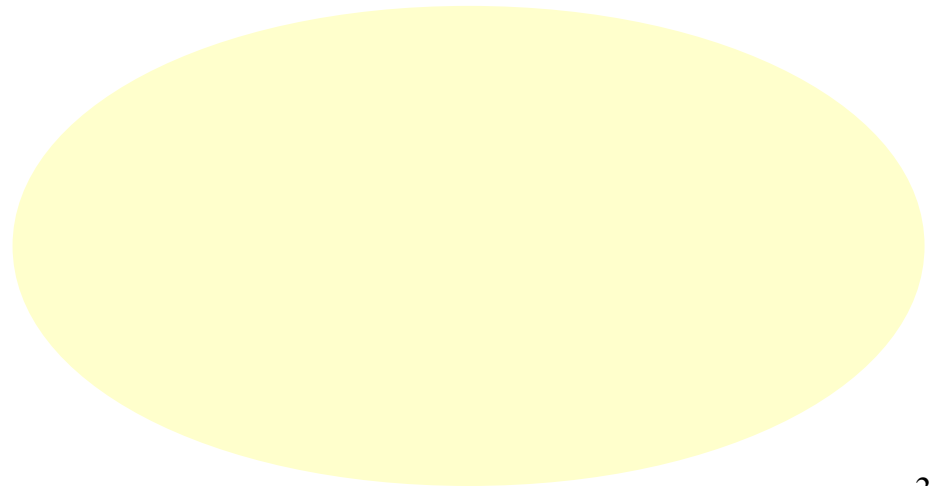
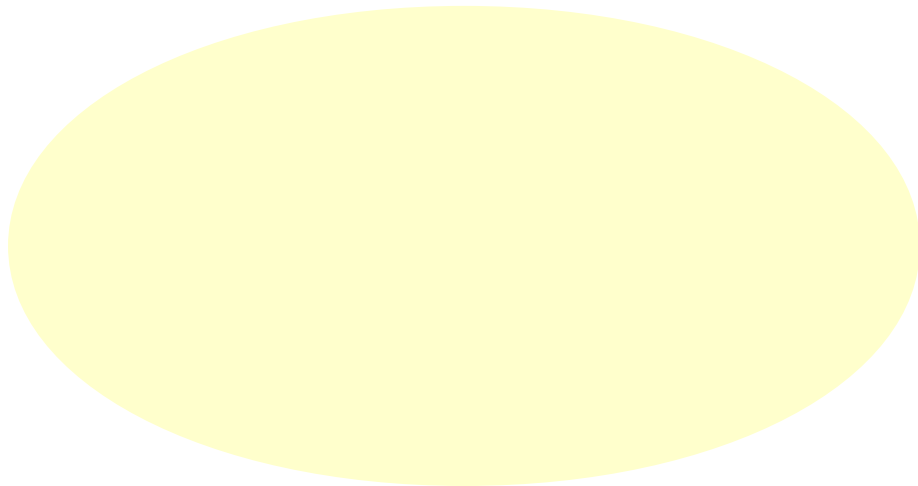
# Regular vs. Irregular Algorithms



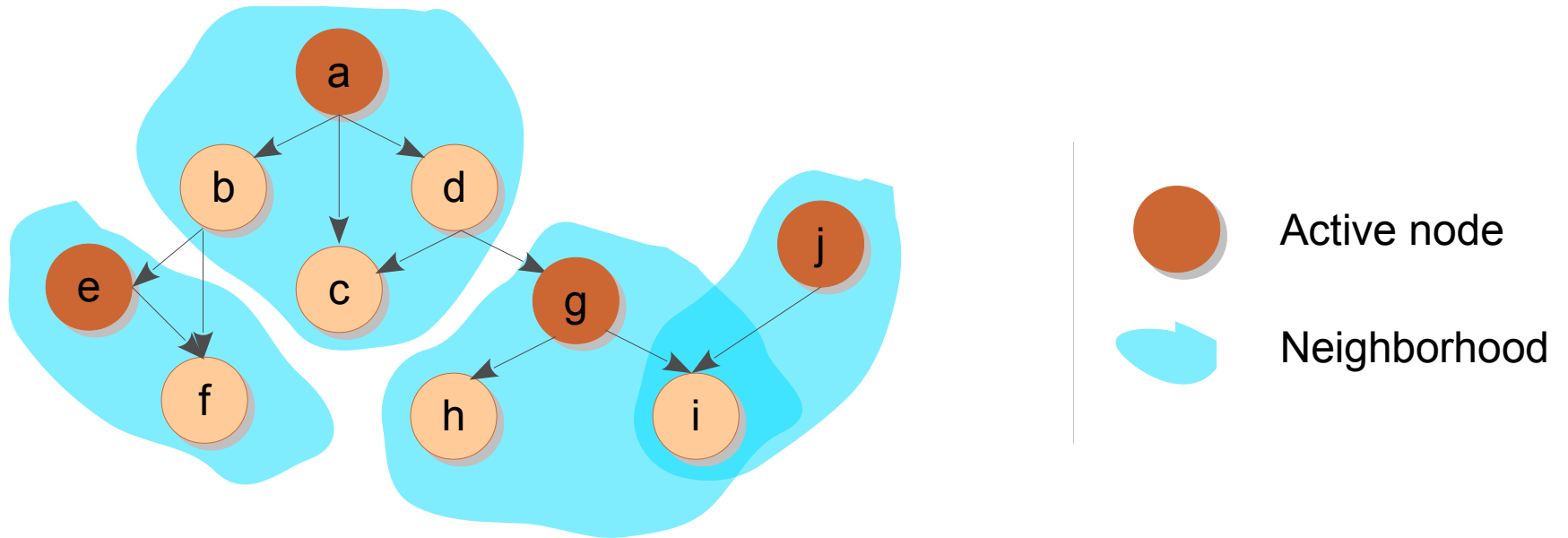
Matrix Multiplication



Shortest Paths Computation



# Dynamic Techniques



Non-overlapping neighborhoods can be processed in parallel.

Overlapping neighborhoods require synchronization.

Leads to optimistic and cautious parallelizations.

# Sequential to Parallel

- Sequential programs often overspecify dependencies.

```
for (int ii = 0; ii < N; ++ii) {  
    process(a[ii]);  
}
```

Processing of  $a[ii + 1]$  is specified after that of  $a[ii]$ .

```
x = y;  
f(a, b);  
while (m < n) {  
    process(m);  
    m = next(m);  
}
```

Processing of assignment, function call and while are sequentially specified.

We need a way to specify that various operations need not be executed in a specific order.

# Unordered Execution

```
for (int ii = 0; ii < N; ++ii) {  
    process(a[ii]);  
}
```



```
forall (e in a) {  
    process(e);  
}
```

```
x = y;  
f(a, b);  
while (m < n) {  
    process(m);  
    m = next(m);  
}
```



```
unordered(  
    x = y;  
    f(a, b);  
    while (m < n) {  
        process(m);  
        m = next(m);  
    }  
);
```

# Sequential to Parallel

- We added unorderedness.
- We added non-determinism.
- We added higher-level information.