

15-745 Optimizing For Data Locality - 1

Seth Copen Goldstein
Seth@cs.cmu.edu
CMU

Based on "A Data Locality Optimizing Algorithm,
Wolf & Lam, PLDI '91

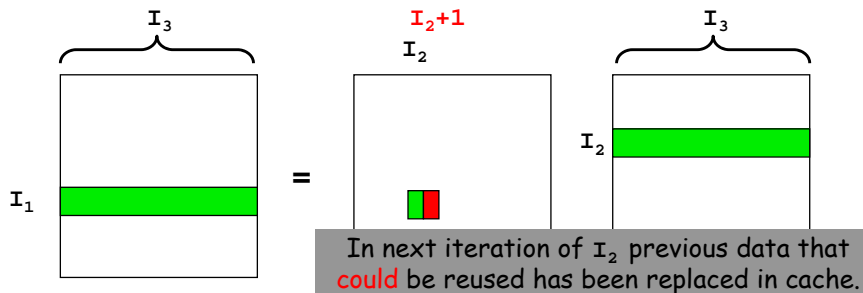
Outline

- The Problem
- Loop Transformations
 - dependence vectors
 - Transformations
 - Unimodular transformations
- Locality Analysis
- SRP

The Issue

- Improve cache reuse in nested loops
- Canonical simple case: Matrix Multiply

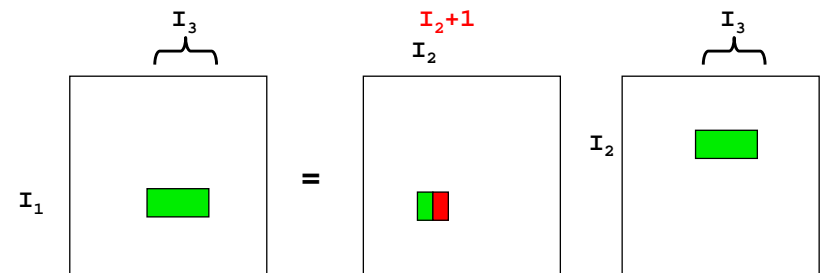
```
for I1 := 1 to n
  for I2 := 1 to n
    for I3 := 1 to n
      C[I1, I3] += A[I1, I2] * B[I2, I3]
```



Tiling solves problem

```
for I1 := 1 to n
  for I2 := 1 to n
    for I3 := 1 to n
      C[I1, I3] += A[I1, I2] * B[I2, I3]
```

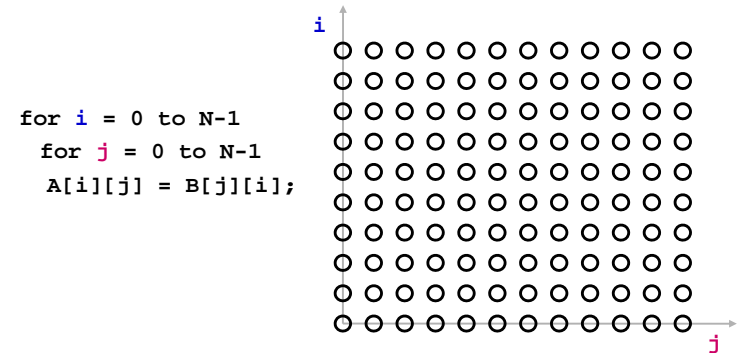
```
for II2 := 1 to n by s
  for II3 := 1 to n by s
    for I1 := 1 to n
      for I2 := II2 to min(II2 + s - 1, n)
        for I3 := II3 to min(II3 + s - 1, n)
          C[I1, I3] += A[I1, I2] * B[I2, I3];
```



The Problem

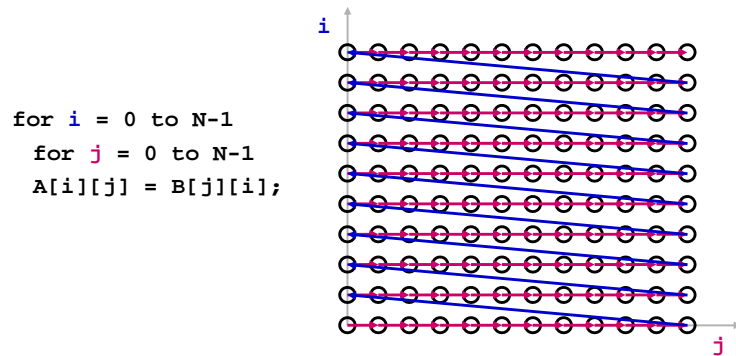
- How to increase locality by transforming loop nest
- Matrix Mult is simple as it is both
 - legal to tile
 - advantageous to tile
- Can we determine the benefit?
(reuse vector space and locality vector space)
- Is it legal (and if so, how) to transform loop?
(unimodular transformations)

Handy Representation: "Iteration Space"



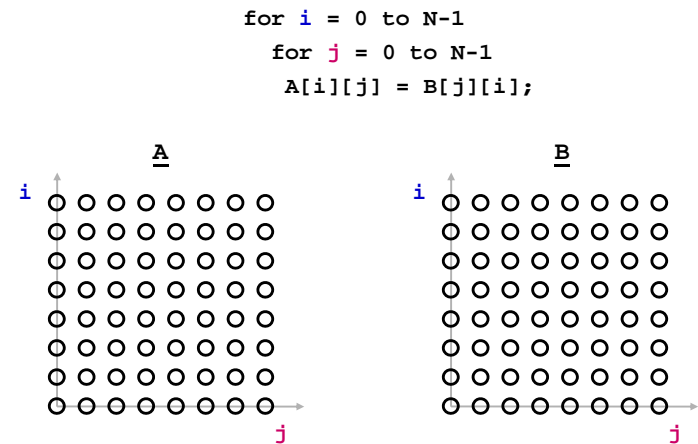
- each position represents an iteration

Visitation Order in Iteration Space



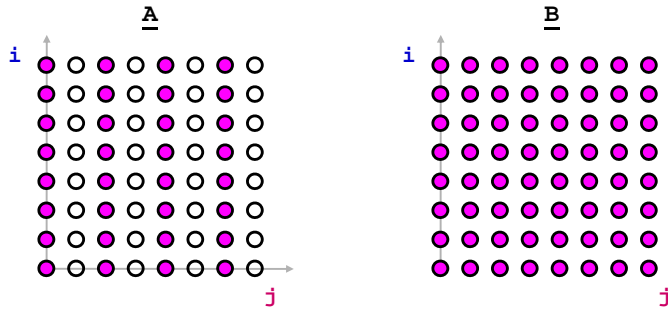
- Note: iteration space is not data space

When Do Cache Misses Occur?



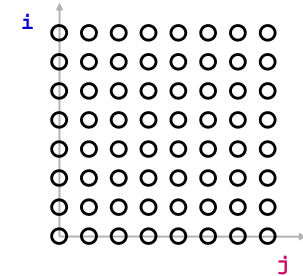
When Do Cache Misses Occur?

```
for i = 0 to N-1
  for j = 0 to N-1
    A[i][j] = B[j][i];
```



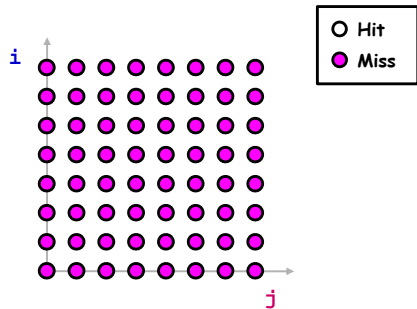
When Do Cache Misses Occur?

```
for i = 0 to N-1
  for j = 0 to N-1
    A[i+j][0] = i*j;
```



When Do Cache Misses Occur?

```
for i = 0 to N-1
  for j = 0 to N-1
    A[i+j][0] = i*j;
```



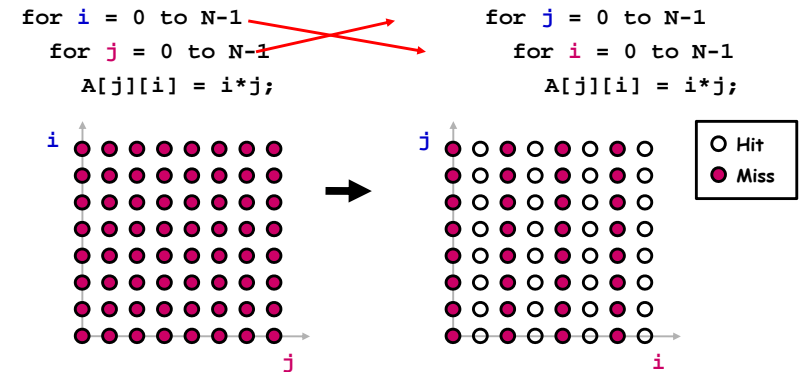
Optimizing the Cache Behavior of Array Accesses

- We need to answer the following questions:
 - when do cache misses occur?
 - use "locality analysis"
 - can we change the order of the iterations (or possibly data layout) to produce better behavior?
 - evaluate the cost of various alternatives
 - does the new ordering/layout still produce correct results?
 - use "dependence analysis"

Examples of Loop Transformations

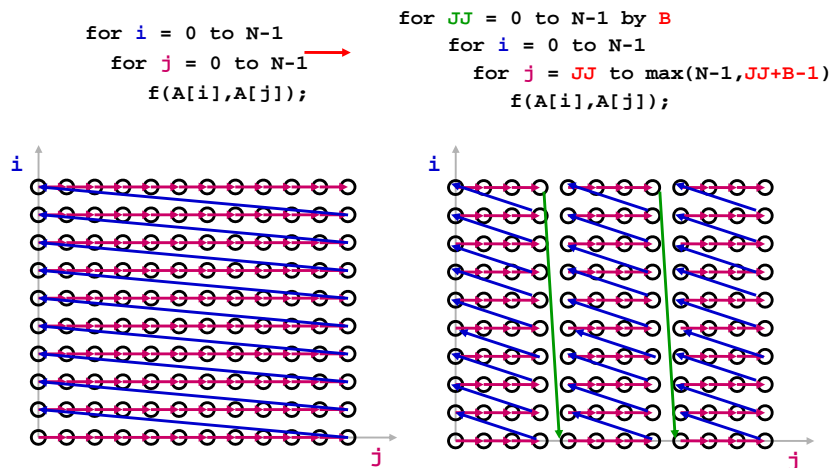
- Loop Interchange
 - Cache Blocking
 - Skewing
 - Loop Reversal
 - ...
- Can improve locality
- Can enable above

Loop Interchange

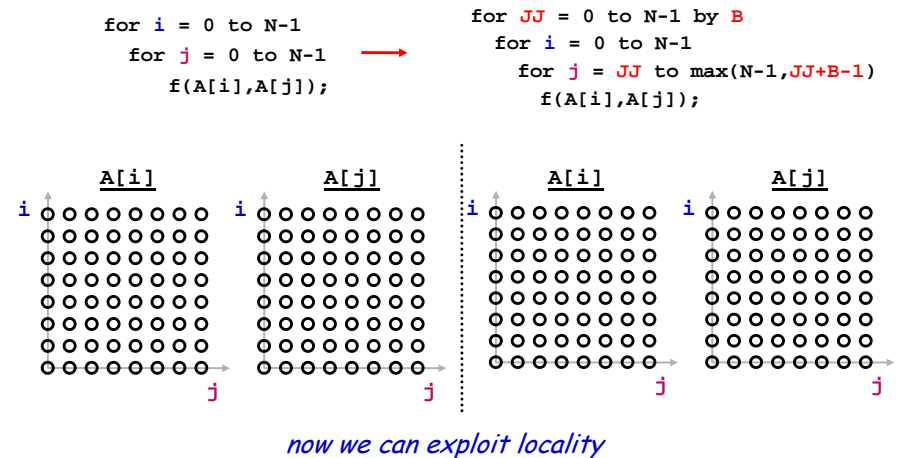


• (assuming N is large relative to cache size)

Impact on Visitation Order in Iteration Space



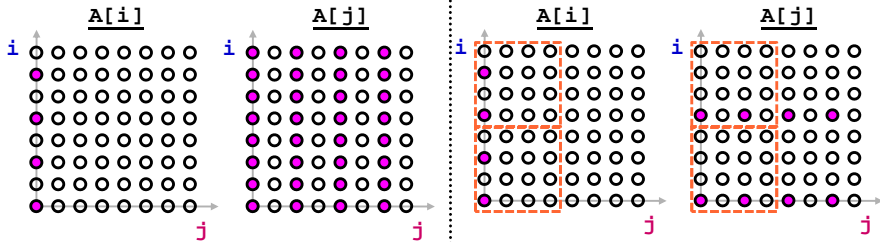
Cache Blocking (aka "Tiling")



Cache Blocking (aka "Tiling")

```
for i = 0 to N-1
  for j = 0 to N-1
    f(A[i],A[j]);
```

```
for JJ = 0 to N-1 by B
  for i = 0 to N-1
    for j = JJ to max(N-1, JJ+B-1)
      f(A[i],A[j]);
```



now we can exploit temporal locality

Cache Blocking in Two Dimensions

```
for i = 0 to N-1
  for j = 0 to N-1
    for k = 0 to N-1
      c[i,k] += a[i,j]*b[j,k];

for JJ = 0 to N-1 by B
  for KK = 0 to N-1 by B
    for i = 0 to N-1
      for j = JJ to max(N-1, JJ+B-1)
        for k = KK to max(N-1, KK+B-1)
          c[i,k] += a[i,j]*b[j,k];
```

- brings square sub-blocks of matrix "b" into the cache
- completely uses them up before moving on

Predicting Cache Behavior through "Locality Analysis"

- Definitions:
 - Reuse: accessing a location that has been accessed in the past
 - Locality: accessing a location that is now found in the cache
- Key Insights
 - Locality only occurs when there is reuse!
 - BUT, reuse does not necessarily result in locality.
 - Why not?

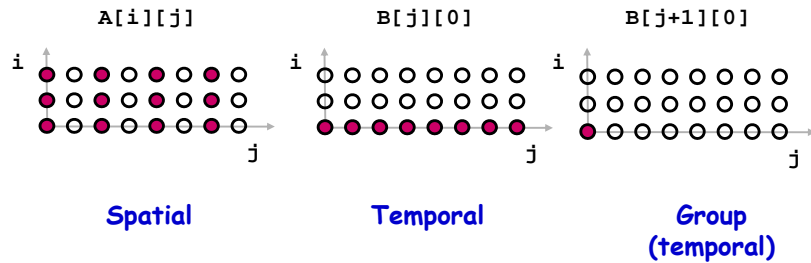
Steps in Locality Analysis

1. Find data reuse
 - if caches were infinitely large, we would be finished
2. Determine "localized iteration space"
 - set of inner loops where the data accessed by an iteration is expected to fit within the cache
3. Find data locality:
 - reuse \supseteq localized iteration space \supseteq locality

Types of Data Reuse/Locality

```
for i = 0 to 2
  for j = 0 to 100
    A[i][j] = B[j][0] + B[j+1][0];
```

○ Hit
● Miss



15-745

© 2005-9 Seth Copen Goldstein

21

Kinds of reuse and the factor

```
for i = 0 to N-1
  for j = 0 to N-1
    f(A[i],A[j]);
```

What kinds of reuse
are there?
 $A[i]$?

$A[j]$?

15-745

© 2005-9 Seth Copen Goldstein

22

Kinds of reuse and the factor

```
for  $I_1 := 0$  to 5
  for  $I_2 := 0$  to 6
     $A[I_2 + 1] = 1/3 * (A[I_2] + A[I_2 + 1] + A[I_2 + 2])$ 
```

15-745

© 2005-9 Seth Copen Goldstein

23

Kinds of reuse and the factor

```
for  $I_1 := 0$  to 5
  for  $I_2 := 0$  to 6
     $A[I_2 + 1] = 1/3 * (A[I_2] + A[I_2 + 1] + A[I_2 + 2])$ 
```

self-temporal in 1, self-spatial in 2
Also, group spatial in 2

What is different about this and previous?

```
for i = 0 to N-1
  for j = 0 to N-1
    f(A[i],A[j]);
```

15-745

© 2005-9 Seth Copen Goldstein

24

Uniformly Generated references

- f and g are indexing functions: $Z^n \rightarrow Z^d$
 - n is depth of loop nest
 - d is dimensions of array, A
- Two references $A[f(i)]$ and $A[g(i)]$ are uniformly generated if

$$f(i) = Hi + c_f \text{ AND } g(i) = Hi + c_g$$

- H is a linear transform
- c_f and c_g are constant vectors

15-745

© 2005-9 Seth Copen Goldstein

25

Eg of Uniformly generated sets

for $I_1 := 0$ to 6
for $I_2 := 0$ to 6
 $A[I_2 + 1] = 1/3 * (A[I_2] + A[I_2 + 1] + A[I_2 + 2])$

$$A[I_2 + 1] \quad [0 \ 1] \begin{bmatrix} I_1 \\ I_2 \end{bmatrix} + [1]$$

$$A[I_2] \quad [0 \ 1] \begin{bmatrix} I_1 \\ I_2 \end{bmatrix} + [0]$$

$$A[I_2 + 2] \quad [0 \ 1] \begin{bmatrix} I_1 \\ I_2 \end{bmatrix} + [2]$$

15-745

© 2005-9 Seth Copen Goldstein

26

Quantifying Reuse

- Why should we quantify reuse?
- How do we quantify locality?

15-745

© 2005-9 Seth Copen Goldstein

27

Quantifying Reuse

- Why should we quantify reuse?
- How do we quantify locality?
- Use vector spaces to identify loops with reuse
- We convert that reuse into locality by making the "best" loop the inner loop
- Metric: memory accesses/iter of innermost loop. No locality \rightarrow mem access

15-745

© 2005-9 Seth Copen Goldstein

28

Self-Temporal

- For a reference, $A[Hi+c]$, there is self-temporal reuse between m and n when $Hm+c=Hn+c$, i.e., $H(r)=0$, where $r=m-n$.
- The direction of reuse is r .
- The self-temporal reuse vector space is: $R_{ST} = \text{Ker } H$
- There is locality if R_{ST} is in the localized vector space.

Recall that for $n \times m$ matrix A ,
the $\text{ker } A = \text{nullspace}(A) = \{x^m \mid Ax = 0\}$

- Reuse is $s^{\dim(R_{ST})}$
- R_{ST} intersect $L = \text{locality}$
- # of mem refs = $1/\text{above}$

Example of self-temporal reuse

```
for I1 := 1 to n
  for I2 := 1 to n
    for I3 := 1 to n
      C[I1, I3] += A[I1, I2] * B[I2, I3]
```

Access	H	ker H	reuse?	Local?
$C[I_1, I_3]$	$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$	$\text{span}\{(0,1,0)\}$	$n \text{ in } I_2$	
$A[I_1, I_2]$	$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}$	$\text{span}\{(0,0,1)\}$		
$B[I_2, I_3]$	$\begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$	$\text{span}\{(1,0,0)\}$		

Self-Spatial

- Occurs when we access in order
 - $A[i,j]$: best gain, l
 - $A[i,j*k]$: best gain, l/k if $|k| \leq l$
- How do we get spatial reuse for UG: H ?

Self-Spatial

- Occurs when we access in order
 - $A[i,j]$: best gain, 1
 - $A[i,j*k]$: best gain, $1/k$ if $|k| \leq 1$
- How do we get spatial reuse for UG: H ?
- Since all but row must be identical, set last row in H to 0, H_s
self-spatial reuse vector space = R_{SS}
 $R_{SS} = \ker H_s$
- Notice, $\ker H \subseteq \ker H_s$
- If, $R_{SS} \cap L = R_{ST} \cap L$, then no additional benefit to SS

15-745

© 2005-9 Seth Copen Goldstein

33

Example of self-spatial reuse

```

for I1 := 1 to n
  for I2 := 1 to n
    for I3 := 1 to n
      C[I1,I3] += A[I1,I2] * B[I2,I3]
    
```

Access	H_s	$\ker H_s$	reuse? Local?
$C[I_1, I_3]$	$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$	$\text{span}\{(0,1,0), (0,0,1)\}$	1/1
$A[I_1, I_2]$	$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$	$\text{span}\{(0,0,1), (0,1,0)\}$	
$B[I_2, I_3]$	$\begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}$	$\text{span}\{(1,0,0), (0,0,1)\}$	

15-745

© 2005-9 Seth Copen Goldstein

34

Self-spatial reuse/locality

- $\text{Dim}(R_{SS})$ is dimensionality of reuse vector space.
- If $R_{SS}=0 \rightarrow$ no reuse
- If $R_{SS}=R_{ST}$ no extra reuse from spatial
- Reuse of each element is $k/s^{\text{dim}(R_{SS})}$ where, s is number of iters per dim.
- $R_{SS} \cap L$ is amount of reuse exploited, therefore number of memory references generated is:
 $k/s^{\text{dim}(R_{ST} \cap L)}$

15-745

© 2005-9 Seth Copen Goldstein

35

Group Temporal

- Two refs $A[Hi+c]$ and $A[Hi+d]$ can have group temporal reuse in L iff
 - they are from same uniformly generated set
 - There is an $r \in L$ s.t. $Hr = c - d$
- if $c-d = r_p$, then there is group temporal reuse, $R_{GT} = \ker H + \text{span}\{r_p\}$
- However, there is no extra benefit if $R_{GT} \cap L = R_{ST} \cap L$

15-745

© 2005-9 Seth Copen Goldstein

36

Example:

```
For i = 1 to n
  for j=i to n
    A[i,j] = 0.2*(A[i,j]+A[i+1,j]+
                 A[i-1,j]+A[i,j+1]+A[i,j-1])
```

If $L = \text{span}\{j\}$, since $\ker H = \emptyset$:

$A[i,j]$ and $A[i,j-1] \rightarrow (0,0)-(0,-1) \in \text{span}\{(0,1)\}$ yes
 $A[i,j-1]$ and $A[i+1,j] \rightarrow (0,-1)-(1,0) \notin \text{span}\{(0,1)\}$ no

Notice equivalence classes

Evaluating group temporal reuse

- Divide all references from a uniformly generated set into equiv classes that satisfy the R_{GT}
- For a particular L and g references
 - Don't count any group reuse when $R_{GT} \cap L = R_{ST} \cap L$
 - number of equiv classes is g_T .
 - Number of mem references is g_T instead of g

Total memory accesses

- For each uniformly generated set localized space, L
line size, z

$$\frac{g_S + (g_T - g_S)/z}{z e_S^{\dim(R_{SS} \cap L)}}$$

$$\text{where } e = \begin{cases} 0 & \text{if } R_{ST} \cap L = R_{SS} \cap L \\ 1 & \text{otherwise} \end{cases}$$

Now what?

- We have a way to characterize
 - Reuse (potential for locality)
 - Local iteration space
- Can we transform loop to take advantage of reuse?
- If so, can we?

Loop Transformation Theory

- Iteration Space
- Dependence vectors
- Unimodular transformations

Loop Nests and the Iter space

- General form of tightly nested loop

```

for I1 := low1 to high1 by step1
  for I2 := low2 to high2 by step2
    ...
    for Ii := lowi to highi by stepi
      ...
      for In := lown to highn by stepn
        Stmts
    
```

- The iteration space is a convex polyhedron in \mathbb{Z}^n bounded by the loop bounds.
- Each iteration is a node in the polyhedron identified by its vector: $\mathbf{p}=(p_1, p_2, \dots, p_n)$

Lexicographical Ordering

- Iterations are executed in lexicographic order.
- for $\mathbf{p}=(p_1, p_2, \dots, p_n)$ and $\mathbf{q}=(q_1, q_2, \dots, q_n)$ if $\mathbf{p} >_k \mathbf{q}$ iff for $1 \leq k \leq n$,

$$\forall 1 \leq i < k, (p_i = q_i) \text{ and } p_k > q_k$$

- For MM:

- (1,1,1), (1,1,2), (1,1,3), ...,
 (1,2,1), (1,2,2), (1,2,3), ...,

...,
 (2,1,1), (2,1,2), (2,1,3), ...

- (1,2,1) $>_2$ (1,1,2), (2,1,1) $>_1$ (1,4,2), etc.

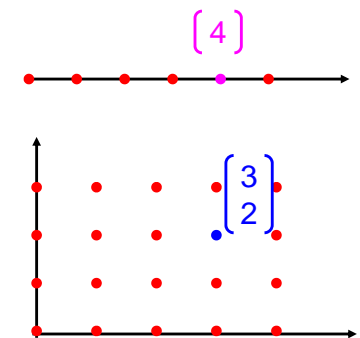
Iteration Space

Every iteration generates a point in an n-dimensional space, where n is the depth of the loop nest.

```

for (i=0; i<n; i++) {
  ...
}
for (i=0; i<n; i++)
  for (j=0; j<4; j++) {
    ...
  }

```



Dependence Vectors

- Dependence vector in an n-nested loop is denoted as a vector: $\mathbf{d}=(d_1, d_2, \dots, d_n)$.
- Each d_i is a possibly infinite range of ints in $[d_i^{\min}, d_i^{\max}]$, where $d_i^{\min} \in \mathbb{Z} \cup \{-\infty\}, d_i^{\max} \in \mathbb{Z} \cup \{\infty\}$ and $d_i^{\min} \leq d_i^{\max}$
- So, a single dep vector represents a set of distance vectors.
- A distance vector defines a distance in the iteration space.
- A dependence vector is a distance vector if each d_i is a singleton.

15-745

© 2005-9 Seth Copen Goldstein

45

Other defs

- Common ranges in dependence vectors
 - $[1, \infty]$ as + or >
 - $[-\infty, -1]$ as - or <
 - $[-\infty, \infty]$ as \pm or *
- A distance vector is the difference between the target and source iterations (for a dependent ref), e.g., $\mathbf{d} = \mathbf{I}_t - \mathbf{I}_s$

15-745

© 2005-9 Seth Copen Goldstein

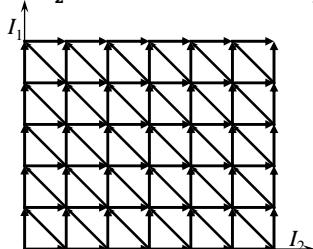
46

Examples

```
for I1 := 1 to n
  for I2 := 1 to n
    for I3 := 1 to n
      C[I1, I3] += A[I1, I2] * B[I2, I3]
```

(0,1,0)

```
for I1 := 0 to 5
  for I2 := 0 to 6
    A[I2 + 1] := 1/3 * (A[I2] + A[I2 + 1] + A[I2 + 2])
```



$D = \{(0,1), (1,0), (1-1)\}$

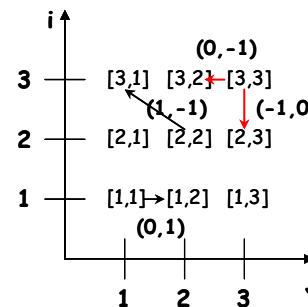
15-745

© 2005-9 Seth Copen Goldstein

47

Plausible Dependence vectors

- A dependence vector is plausible iff it is lexicographically non-negative.
- All sequential programs have plausible dependence vectors. Why?
- Plausible: (1,-1)
- implausible (-1,0)



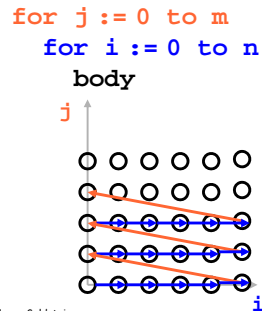
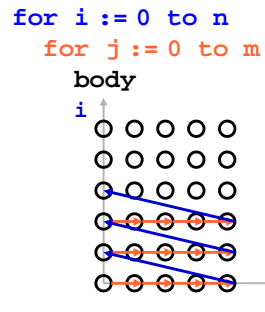
15-745

© 2005-9 Seth Copen Goldstein

48

Loop Transforms

- A loop transformation changes the order in which iterations in the iteration space are visited.
- For example, Loop Interchange



15-745

© 2005-9 Seth Copen Goldstein

49

Unimodular Transforms

- Interchange
permute nesting order
- Reversal
reverse order of iterations
- Skewing
scale iterations by an outer loop index

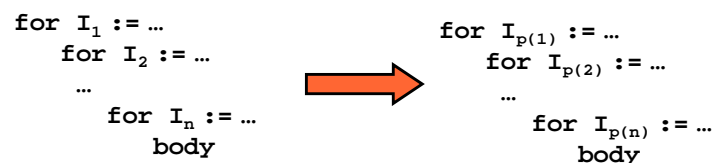
15-745

© 2005-9 Seth Copen Goldstein

50

Interchange

- Change order of loops
- For some permutation p of $1 \dots n$



- When is this legal?

15-745

© 2005-9 Seth Copen Goldstein

51

Transform and matrix notation

- If dependences are vectors in iter space, then transforms can be represented as matrix transforms
- E.g., for a 2-deep loop, interchange is:

$$T = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} p_1 \\ p_2 \end{bmatrix} = \begin{bmatrix} p_2 \\ p_1 \end{bmatrix}$$

- Since, T is a linear transform, Td is transformed dependence:

$$\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \end{bmatrix} = \begin{bmatrix} d_2 \\ d_1 \end{bmatrix}$$

15-745

© 2005-9 Seth Copen Goldstein

52

Reversal

- Reversal of i^{th} loop reverses its traversal, so it can be represented as:

15-745

© 2005-9 Seth Copen Goldstein

53

Reversal

- Reversal of i^{th} loop reverses its traversal, so it can be represented as: Diagonal matrix with i^{th} element = -1.
- For 2 deep loop, reversal of outermost is:

$$T = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \quad \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} p_1 \\ p_2 \end{bmatrix} = \begin{bmatrix} -p_1 \\ p_2 \end{bmatrix}$$

15-745

© 2005-9 Seth Copen Goldstein

54

Skewing

- Skew loop I_j by a factor f w.r.t. loop I_i maps

$$(p_1, \dots, p_i, \dots, p_j, \dots) \quad (p_1, \dots, p_i, \dots, p_j + fp_i, \dots)$$

- Example for 2D

$$T = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \quad \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} p_1 \\ p_2 \end{bmatrix} = \begin{bmatrix} p_1 \\ p_2 + p_1 \end{bmatrix}$$

15-745

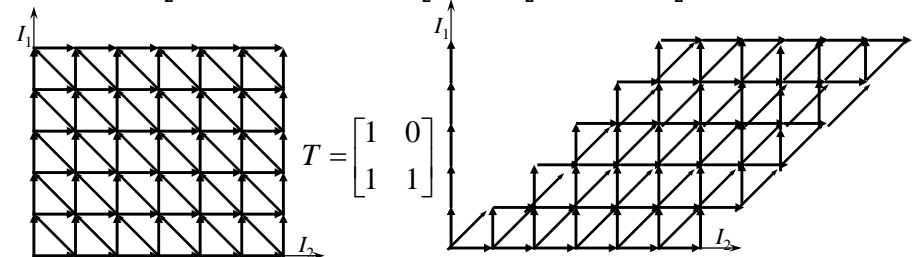
© 2005-9 Seth Copen Goldstein

55

Loop Skewing Example

```
for I1 := 0 to 5
  for I2 := 0 to 6
    A[I2 + 1] := 1/3 * (A[I2] + A[I2 + 1] + A[I2 + 2])
```

$D = \{(0,1), (1,0), (1,-1)\}$



```
for I1 := 0 to 5
  for I2 := I1 to 6+I1
    A[I2-I1+1] := 1/3 * (A[I2-I1] + A[I2-I1+1] + A[I2-I1+2])
```

$D = \{(0,1), (1,1), (1,0)\}$

15-745

© 2005-9 Seth Copen Goldstein

56

But...is the transform legal?

- Distance/direction vectors give a partial order among points in the iteration space
- A loop transform changes the order in which 'points' are visited
- The new visit order must respect the dependence

15-745

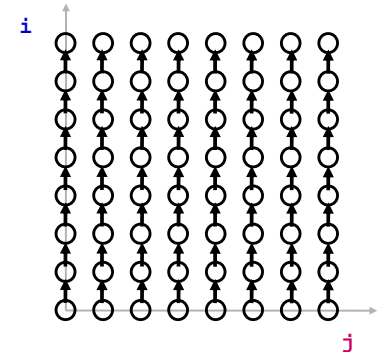
© 2005-9 Seth Copen Goldstein

57

But...is the transform legal?

- Loop reversal ok?
- Loop interchange ok?

```
for i = 0 to N-1
  for j = 0 to N-1
    A[i+1][j] += A[i][j];
```



15-745

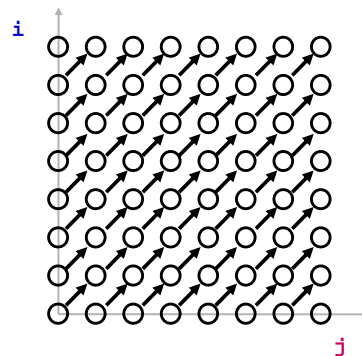
© 2005-9 Seth Copen Goldstein

58

But...is the transform legal?

- Loop reversal ok?
- Loop interchange ok?

```
for i = 0 to N-1
  for j = 0 to N-1
    A[i+1][j+1] += A[i][j];
```



15-745

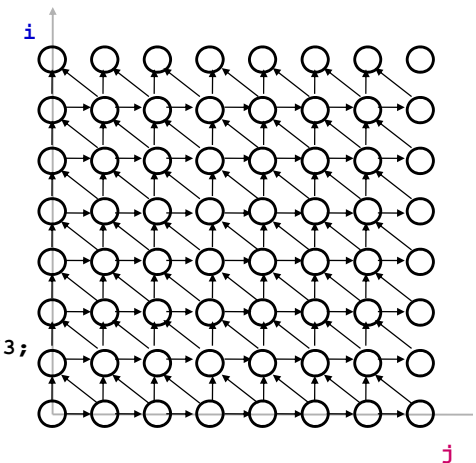
© 2005-9 Seth Copen Goldstein

59

But...is the transform legal?

- What other visit order is legal here?

```
for i = 0 to TS
  for j = 0 to N-2
    A[j+1] =
      (A[j] + A[j+1] + A[j+2])/3;
```



15-745

© 2005-9 Seth Copen Goldstein

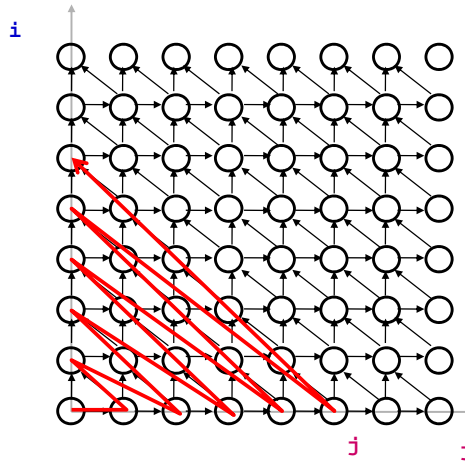
60

But...is the transform legal?

- What other visit order is legal here?

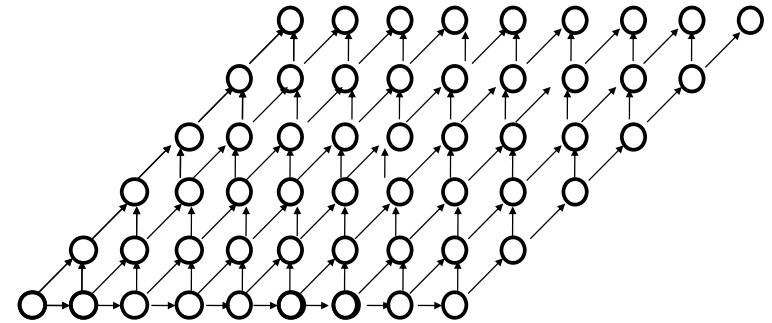
```

for i = 0 to TS
  for j = 0 to N-2
    A[j+1] =
      (A[j] + A[j+1] + A[j+2])/3;
  
```



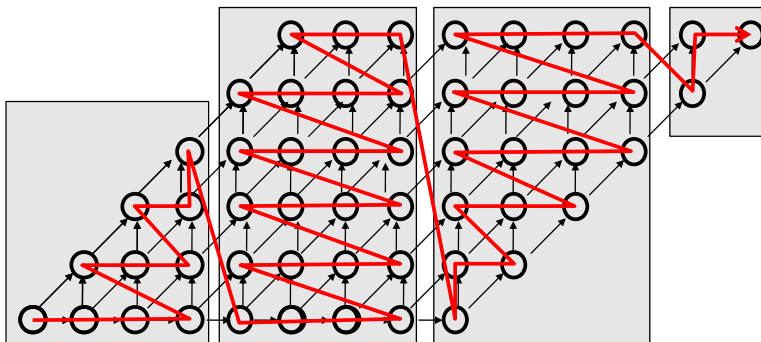
But...is the transform legal?

- Skewing...



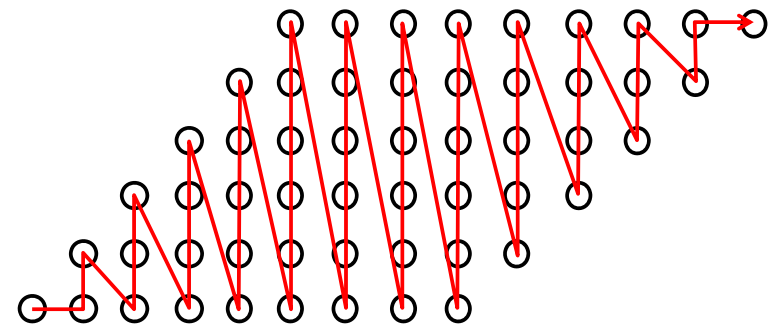
But...is the transform legal?

- Skewing...now we can block



But...is the transform legal?

- Skewing...now we can loop interchange



Unimodular transformations

- Express loop transformation as a matrix multiplication
- Check if any dependence is violated by multiplying the distance vector by the matrix - if the resulting vector is still lexicographically positive, then the involved iterations are visited in an order that respects the dependence.

Reversal

$$\begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$$

Interchange

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

Skew

$$\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$$

"A Data Locality Optimizing Algorithm", M.E.Wolf and M.Lam

Next Time

- Putting it all together: SRP
- Other loop transformations for locality

Linear Algebra

- Vector Spaces
- Linear Combinations
- dimensions
- Spans
- Kernels

Vector Spaces

- \mathbf{n} is a point in n -space
- $V = \{ \mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m \}$ is a finite set of n -vectors over $m \mathbb{R}^n$.
- Linear combination of vectors of V is a vector \mathbf{x} as defined by
$$\mathbf{x} = \alpha_1 \mathbf{v}_1 + \alpha_2 \mathbf{v}_2 + \dots + \alpha_m \mathbf{v}_m$$
where α_i are real numbers.
- V is linearly dependent if a combination results in the $\mathbf{0}$ vector, otherwise it is linearly independent.

Dim and Basis

- dimensionality of V is $\dim(V)$
the number of independent vectors in V
- A basis for an m -dimensional vector space is a set of linearly independent vectors such that every point in V can be expressed as a linear comb of the vectors in the basis.
 - The vectors in the basis are called basis vectors

Subspaces and span

- Let V be a set of vectors
- The subspace spanned by V , $\text{span}(V)$, is a subset of \mathfrak{R}^n such that
 - $V \subseteq \text{span}(V)$
 - $\mathbf{x}, \mathbf{y} \in \text{span}(V) \Rightarrow \mathbf{x} + \mathbf{y} \in \text{span}(V)$
 - $\mathbf{x} \in \text{span}(V)$ and $\alpha \in \mathfrak{R} \Rightarrow \alpha \mathbf{x} \in \text{span}(V)$

Range, Span, Kernel

- A matrix A can be viewed as a set of column vectors.
- Range $A^{n \times m}$ is $\{A\mathbf{x} \mid \mathbf{x} \in \mathfrak{R}^m\}$
- $\text{span}(A) = \text{Range } A^{n \times m}$
- $\text{nullspace}(A) = \ker(A) = \ker(A^{n \times m}) = \{\mathbf{x}^m \mid A\mathbf{x} \in \mathbf{0}\}$
- $\text{rank}(A) = \dim(\text{span}(A))$
- $\text{nullity}(A) = \dim(\ker(A))$
- $\text{rank}(A) + \text{nullity}(A) = n$, for $A^{n \times m}$