



Optimization and Parallelization of Sequential Programs

Lecture 7

Christoph Kessler
IDA / PELAB
Linköping University
Sweden

Outline



Towards (semi-)automatic parallelization of sequential programs

- Data dependence analysis for loops
- Some loop transformations
 - Loop invariant code hoisting, loop unrolling, loop fusion, loop interchange, loop blocking and tiling
- Static loop parallelization
- Run-time loop parallelization
 - Doacross parallelization, Inspector-executor method
- Speculative parallelization (as time permits)
- Auto-tuning (later, if time)

Foundations: Control and Data Dependence



- Consider statements S, T in a sequential program ($S=T$ possible)
 - Scope of analysis is typically a function, i.e. intra-procedural analysis
 - Assume that a control flow path $S \dots T$ is possible
 - Can be done at arbitrary granularity (instructions, operations, statements, compound statements, program regions)
 - Relevant are only the read and write effects on memory (i.e. on program variables) by each operation, and the effect on control flow

- **Control dependence** $S \rightarrow T$, if the fact whether T is executed may depend on S (e.g. condition)

- Implies that relative execution order $S \rightarrow T$ must be preserved when restructuring the program
- Mostly obvious from nesting structure in well-structured programs, but more tricky in arbitrary branching code (e.g. assembler code)

Example:

```
S: if (...) {
...
T: ...
...
}
```

Foundations: Control and Data Dependence



- **Data dependence** $S \rightarrow T$, if statement S may execute (dynamically) before T and both may access the same memory location and at least one of these accesses is a write

Example:

```
S: z = ... ;
...
T: ... = ..z.. ;
```

(flow dependence)

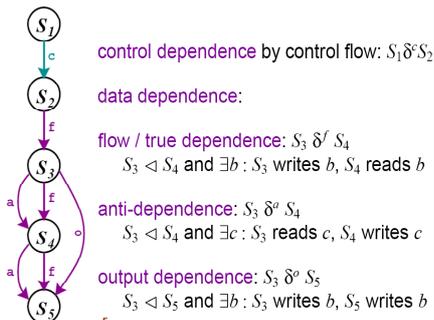
- Means that execution order " S before T " must be preserved when restructuring the program
- In general, only a conservative over-estimation can be determined statically
- **flow dependence:** (RAW, read-after-write)
 - ▶ S may write a location z that T may read
- **anti dependence:** (WAR, write-after-read)
 - ▶ S may read a location x that T may overwrites
- **output dependence:** (WAW, write-after-write)
 - ▶ both S and T may write the same location

Dependence Graph



- **(Data, Control, Program) Dependence Graph:** Directed graph, consisting of all statements as vertices and all (data, control, any) dependences as edges.

```
S1: if (e) goto S3
S2: a ← ...
S3: b ← a * c
S4: c ← b * f
S5: b ← x + f
```



Why Loop Optimization and Parallelization



Loops are a promising object for program optimizations, including automatic parallelization:

- High execution frequency
 - Most computation done in (inner) loops
 - Even small optimizations can have large impact (cf. Amdahl's Law)
- Regular, repetitive behavior
 - compact description
 - *relatively* simple to analyze statically
- Well researched

Loop Optimizations – General Issues



- Move loop invariant computations out of loops
- Modify the order of iterations or parts thereof

Goals:

- Improve data access locality
- Faster execution
- Reduce loop control overhead
- Enhance possibilities for loop parallelization or vectorization

Only transformations that preserve the program semantics (its input/output behavior) are admissible

- Conservative (static) criterium: preserve data dependences
- Need data dependence analysis for loops



Data Dependence Analysis for Loops

A more formal introduction

Data Dependence Analysis – Overview



- Important for loop optimizations, vectorization and parallelization, instruction scheduling, data cache optimizations
- Conservative approximations to disjointness of pairs of memory accesses
 - weaker than data-flow analysis
 - but generalizes nicely to the level of individual array element
- Loops, loop nests
 - Iteration space
 - Array subscripts in loops
 - Index space
- Dependence testing methods
- Data dependence graph
- Data + control dependence graph
 - Program dependence graph

Precedence relation between statements



S_1 statically (textually) precedes S_2 $S_1 \text{ pred } S_2$

S_1 dynamically precedes S_2 $S_1 < S_2$

Within loops, loop nests: $\text{pred} \neq <$

```

 $S_1$ :  $s \leftarrow 0$ 
  for  $i$  from 1 to  $n$  do
 $S_2$ :    $s \leftarrow s + a[i]$ 
 $S_3$ :    $a[i] \leftarrow s$ 
  od
    
```

Data Dependence Graph

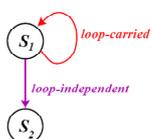


- Data dependence graph for straight-line code ("basic block", no branching) is always *acyclic*, because relative execution order of statements is forward only.
- Data dependence graph for a loop:
 - Dependence edge $S \rightarrow T$ if a dependence *may* exist for some pair of instances (iterations) of S, T
 - Cycles possible
 - Loop-independent versus loop-carried dependences

Example:

```

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



(assuming we know statically that arrays a and b do not intersect)

Loop Iteration Space

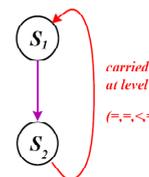


Beyond basic blocks: $\text{pred} \neq <$

Canonical loop nest: (HIR code)

```

for  $i_1$  from 1 to  $n_1$  do
  for  $i_2$  from 1 to  $n_2$  do
    ...
    for  $i_k$  from 1 to  $n_k$  do
       $S_1(i_1, \dots, i_k)$ :  $A[i_1, 2 * i_3] \leftarrow B[i_2, i_3] + 1$ 
       $S_2(i_1, \dots, i_k)$ :  $B[i_2, i_3 + i_4] \leftarrow 2 * A[i_1, 2 * i_3]$ 
    od
  od
od
    
```



Iteration space: $IS = [1..n_1] \times [1..n_2] \times \dots \times [1..n_k]$

(the simplest case: rectangular, static loop bounds)

Iteration vector $\vec{i} = \langle i_1, \dots, i_k \rangle \in IS$



Example

```

for i from 2 to 9 do
S1 X[i] ← Y[i] + Z[i]
S2 A[i] ← X[i-1] + 1
od

```

(assuming that we statically know that arrays A, X, Y, Z do not intersect, otherwise there might be further dependences)

	$i = 2$	$i = 3$	$i = 4$...
S_1	$X[2] \leftarrow Y[2] + Z[2]$	$X[3] \leftarrow Y[3] + Z[3]$	$X[4] \leftarrow Y[4] + Z[4]$...
S_2	$A[2] \leftarrow X[1] + 1$	$A[3] \leftarrow X[2] + 1$	$A[4] \leftarrow X[3] + 1$...

There is a loop-carried, forward, flow dependence from S_1 to S_2 .



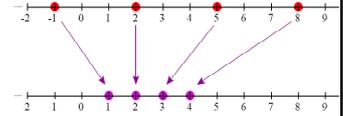
Loop Normalization

Given a loop of the form

```

for I from L to U step S do
... I ...
od

```



normalize the loop:

- lower bound 0 (C) resp. 1 (Fortran)
- step size +1

→ update all occurrences of the loop counter I by $i * S - S + L$

```

for i from 1 to (U-L+S)/S step 1 do
... (i*S-S+L) ...
od
I ← i*S-S+L

```



Dependence Distance and Direction

Lexicographic order on iteration vectors → dynamic execution order:

$$S_1((i_1, \dots, i_k)) \prec S_2((j_1, \dots, j_k)) \text{ iff}$$

either $S_1 \text{ pred } S_2$ and $\langle i_1, \dots, i_k \rangle \leq_{lex} \langle j_1, \dots, j_k \rangle$
or $S_1 = S_2$ and $\langle i_1, \dots, i_k \rangle <_{lex} \langle j_1, \dots, j_k \rangle$

distance vector $\vec{d} = \vec{j} - \vec{i} = \langle j_1 - i_1, \dots, j_k - i_k \rangle$

direction vector $dirv = \text{sgn}(\vec{j} - \vec{i}) = \langle \text{sgn}(j_1 - i_1), \dots, \text{sgn}(j_k - i_k) \rangle$
in terms of symbols = $\langle \langle \rangle \leq \rangle^*$

Example: $S_1((i_1, i_2, i_3, i_4)) \delta^f S_2((i_1, i_2, i_3, i_4))$
distance vector $\vec{d} = \langle 0, 0, 0, 0 \rangle$, direction vector $dirv = \langle -, -, -, - \rangle$,
loop-independent dependence

Example: $S_2((i_1, i_2, i_3, i_4)) \delta^f S_1((i_1, i_2, i_3, i_4))$
distance vector $\vec{d} = \langle 0, 0, ?, 0 \rangle$, direction vector $dirv = \langle =, =, >, = \rangle$,
loop-carried dependence (carried by i_3 -loop / at level 3)



Dependence Equation System

One-dimensional array A accessed in k nested loops: $S_1 : \dots A[f(\vec{i})] \dots$
 $S_2 : \dots A[g(\vec{j})] \dots$

Is there a dependence between $S_1(\vec{i})$ and $S_2(\vec{j})$ for some $\vec{i}, \vec{j} \in ItS$?

typically f, g linear: $f(\vec{i}) = a_0 + \sum_{i=1}^k a_i i_i$, $g(\vec{j}) = b_0 + \sum_{j=1}^k b_j j_j$

Exist $\vec{i}, \vec{j} \in \mathbb{Z}^k$ with $f(\vec{i}) = g(\vec{j})$, i.e., $a_0 + \sum_{i=1}^k a_i i_i = b_0 + \sum_{j=1}^k b_j j_j$, dep. equation

subject to $\vec{i}, \vec{j} \in ItS$, i.e.,

$$\begin{aligned}
1 \leq i_1 \leq n_1, & \quad 1 \leq j_1 \leq m_1, \\
& \quad \vdots \\
1 \leq i_k \leq n_k, & \quad 1 \leq j_k \leq m_k
\end{aligned}$$

iter. space constraints: linear inequalities

⇒ constrained linear Diophantine equation system → ILP (NP-complete)



Linear Diophantine Equations

$$\sum_{j=1}^n a_j x_j = c$$

where $n \geq 1$, $c, a_j \in \mathbb{Z}$, $\exists j : a_j \neq 0$, $x_i \in \mathbb{Z}$

Example 1: $x + 4y = 1$
has infinitely many solutions, e.g. $x = 5$ and $y = -1$.

Example 2: $5x - 10y = 2$
has no solution in \mathbb{Z} : absolute term must be multiple of 5

Theorem:
 $\sum_{j=1}^n a_j x_j = c$ has a solution iff $\text{gcd}(a_1, \dots, a_n) | c$.

Proof: see e.g. [Zima/Chapman p. 143]



Dependence Testing, 1: GCD-Test

Often, a simple test is sufficient to prove independence: e.g.,

gcd-test [Banerjee'76], [Towle'76]:

independence if $\text{gcd} \left(\bigcup_{i=1}^n \{a_i, b_i\} \right) \nmid \sum_{i=0}^n (a_i - b_i)$

constraints on ItS not considered

Example: for i from 1 to 4 do

```

S1 : b[i] ← a[3*i-5] + 2
S2 : a[2*i+1] ← 1.0/i

```

solution to $2i + 1 = 3j - 5$ exists in \mathbb{Z} as $\text{gcd}(3, 2) | (-5 - 1 + 3 - 2)$
not checked whether such i, j exist in $\{1, \dots, 4\}$

For multidimensional arrays?



subscript-wise test vs. linearized indexing

for $i \dots$ for $i \dots$

$S_1 : \dots A[x[i], 2 * i] \dots$ $S_1 : \dots A[i, i] \dots$ $A[i * (s_1 + 1)]$

$S_2 : \dots A[y[i], 2 * i + 1] \dots$ $S_2 : \dots A[i, i + 1] \dots$ $A[i * (s_1 + 1) + 1]$

Moreover:

Hierarchical structuring of dependence tests [Burke/Cytron'86]

Survey of Dependence Tests



gcd test

separability test (gcd test for special case, exact)

Banerjee-Wolfe test [Banerjee'88] rational solution in ITS

Delta-test [Goff/Kennedy/Tseng'91]

Power test [Wolfe/Tseng'91]

Simple Loop Residue test [Maydan/Hennessy/Lam'91]

Fourier-Motzkin Elimination [Maydan/Hennessy/Lam'91]

Omega test [Pugh/Wonnacott'92]



Loop Transformations and Parallelization

Loop Optimizations – General Issues



- Move loop invariant computations out of loops
- Modify the order of iterations or parts thereof

Goals:

- Improve data access locality
- Faster execution
- Reduce loop control overhead
- Enhance possibilities for loop parallelization or vectorization

Only transformations that preserve the program semantics (its input/output behavior) are admissible

- Conservative (static) criterium: preserve data dependences
- Need data dependence analysis for loops

Some important loop transformations



- Loop normalization
- Loop parallelization
- Loop invariant code hoisting
- Loop interchange
- Loop fusion vs. Loop distribution / fission
- Strip-mining / loop tiling / blocking vs. Loop linearization
- Loop unrolling, unroll-and-jam
- Loop peeling
- Index set splitting, Loop unswitching
- Scalar replacement, Scalar expansion
- Later: Software pipelining
- More: Cycle shrinking, Loop skewing, ...

Loop Invariant Code Hoisting



- Move loop invariant code out of the loop

- Compilers can do this automatically *if* they can statically find out what code is loop invariant

- Example:

<pre>for (i=0; i<10; i++) a[i] = b[i] + c / d;</pre>		<pre>tmp = c / d; for (i=0; i<10; i++) a[i] = b[i] + tmp;</pre>
---	--	--

Loop Unrolling

Loop unrolling

- Can be enforced with compiler options e.g. `-funroll=2`

Example:

```

for (i=0; i<50; i++) {
    a[i] = b[i];
}
    
```

Unroll by 2:

```

for (i=0; i<50; i+=2) {
    a[i] = b[i];
    a[i+1] = b[i+1];
}
    
```

- Reduces loop overhead (total # comparisons, branches, increments)
- Longer loop body may enable further local optimizations (e.g. common subexpression elimination, register allocation, instruction scheduling, using SIMD instructions)
- longer code

→ Exercise: Formulate the unrolling rule for statically unknown upper loop limit

Loop Interchange (1)

- For properly nested loops (statements in innermost loop body only)

Example 1:

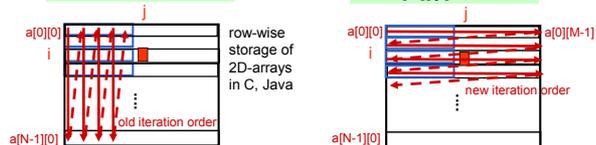
```

for (j=0; j<M; j++)
    for (i=0; i<N; i++)
        a[i][j] = 0.0;
    
```

→

```

for (i=0; i<N; i++)
    for (j=0; j<M; j++)
        a[i][j] = 0.0;
    
```



- Can improve data access locality in memory hierarchy (fewer cache misses / page faults)

C. Kessler, IDA, Linköping universitet.

26

Foundations: Loop-Carried Data Dependences

- Recall: Data dependence $S \rightarrow T$, if operation S may execute (dynamically) before operation T and both may access the same memory location and at least one of these accesses is a write

```

S: z = ...;
T: ... = ..z..;
    
```

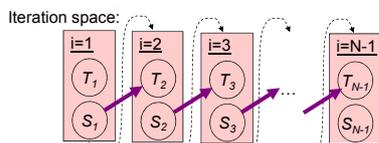
- In general, only a conservative over-estimation can be determined statically.

- Data dependence $S \rightarrow T$ is called **loop carried** by a loop L if the data dependence $S \rightarrow T$ may exist for instances of S and T in different iterations of L .

Example:

```

L: for (i=1; i<N; i++) {
    Ti: ... = x[i-1];
    Si: x[i] = ...;
}
    
```



→ partial order between the operation instances resp. iterations

C. Kessler, IDA, Linköping universitet.

27

Loop Interchange (2)

- Be careful with loop carried data dependences!

Example 2:

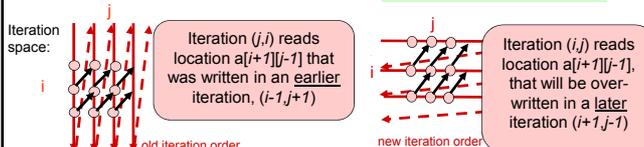
```

for (j=1; j<M; j++)
    for (i=0; i<N; i++)
        a[i][j] = ...a[i+1][j-1]...;
    
```

✗

```

for (i=0; i<N; i++)
    for (j=1; j<M; j++)
        a[i][j] = ...a[i+1][j-1]...;
    
```



- Interchanging the loop headers would violate the partial iteration order given by the data dependences

C. Kessler, IDA, Linköping universitet.

28

Loop Interchange (3)

- Be careful with loop-carried data dependences!

Example 3:

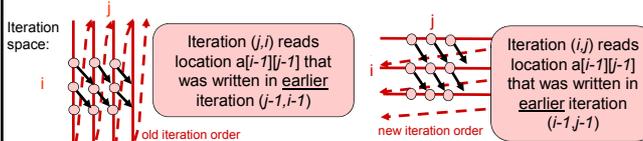
```

for (j=1; j<M; j++)
    for (i=1; i<N; i++)
        a[i][j] = ...a[i-1][j-1]...;
    
```

OK →

```

for (i=1; i<N; i++)
    for (j=1; j<M; j++)
        a[i][j] = ...a[i-1][j-1]...;
    
```



- Generally: Interchanging loop headers is only admissible if loop-carried dependences have the same direction for all loops in the loop nest (all directed along or all against the iteration order)

C. Kessler, IDA, Linköping universitet.

29

Loop Fusion

- Merge subsequent loops with same header

- Safe if neither loop carries a (backward) dependence

Example:

```

for (i=0; i<N; i++)
    a[i] = ...;
for (i=0; i<N; i++)
    ... = ... a[i] ...;
    
```

→

```

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

For N sufficiently large, $a[i]$ will no longer be in the cache at this time

OK – Read of $a[i]$ still after write of $a[i]$, for all i

- Can improve data access locality and reduces number of branches

C. Kessler, IDA, Linköping universitet.

30

Loop Iteration Reordering



A transformation that reorders the iterations of a level- k -loop, without making any other changes, is valid if the loop carries no dependence.

Example:

```
for (i=1; i<n; i++)
  for (j=1; j<m; j++)
    for (k=1; k<r; k++)
S:   a[i][j][k] = ... a[i][j-1][k] ...      (=, <, =)
```

j-loop carries a dependence, its iteration order must be preserved

Loop Parallelization



A transformation that reorders the iterations of a level- k -loop, without making any other changes, is valid if the loop carries no dependence.

Example:

```
for (i=1; i<n; i++)
  for (j=1; j<m; j++)
    for (k=1; k<r; k++)
S:   a[i][j][k] = ... a[i][j-1][k] ...      (=, <, =)
```

j-loop carries a dependence, its iteration order must be preserved

It is valid to convert a sequential loop to a parallel loop if it does not carry a dependence.

Example:

```
for (i=1; i<n; i++)
S:   b[i] = 2 * c[i];
```

Loop parallelization →

```
forall ( i, 1, n, p )
  b[i] = 2 * c[i];
```

Remark on Loop Parallelization



Introducing temporary copies of arrays can remove some antidependences to enable automatic loop parallelization

Example:

```
for (i=0; i<n; i++)
  a[i] = a[i] + a[i+1];
```

The loop-carried dependence can be eliminated:

```
for (i=0; i<n; i++)
  aold[i+1] = a[i+1];
for (i=0; i<n; i++)
  a[i] = a[i] + aold[i+1];
```



Strip Mining / Loop Blocking / -Tiling



```
for (i=0; i<n; i++)
  a[i] = b[i] + c[i];
```

↓ loop blocking with block size s

```
for (i1=0; i1<n; i1+=s) // loop over blocks
  for (i2=0; i2<min(n-i1,s); i2++) // loop within blocks
    a[i1+i2] = b[i1+i2] + c[i1+i2];
```

Tiling = blocking in multiple dimensions + loop interchange

Goal: increase locality; support vectorization (vector registers)

Reverse transformation: Loop linearization

Tiled Matrix-Matrix Multiplication (1)

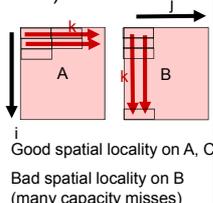


Matrix-Matrix multiplication $C = A \times B$ here for square ($n \times n$) matrices C, A, B , with n large ($\sim 10^3$):

$$C_{ij} = \sum_{k=1..n} A_{ik} B_{kj} \quad \text{for all } i, j = 1..n$$

Standard algorithm for Matrix-Matrix multiplication (here without the initialization of C -entries to 0):

```
for (i=0; i<n; i++)
  for (j=0; j<n; j++)
    for (k=0; k<n; k++)
      C[i][j] += A[i][k] * B[k][j];
```



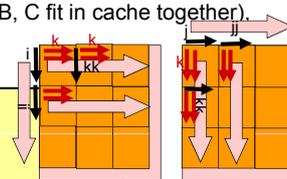
Tiled Matrix-Matrix Multiplication (2)



Block each loop by block size S (choose S so that a block of A, B, C fit in cache together), then interchange loops

Code after tiling:

```
for (ii=0; ii<n; ii+=S)
  for (jj=0; jj<n; jj+=S)
    for (kk=0; kk<n; kk+=S)
      for (i=ii; i < ii+S; i++)
        for (j=jj; j < jj+S; j++)
          for (k=kk; k < kk+S; k++)
            C[i][j] += A[i][k] * B[k][j];
```



Good spatial locality for A, B and C

Remark on Locality Transformations



- An alternative can be to change the data layout rather than the control structure of the program
 - **Example:** Store matrix B in transposed form, or, if necessary, consider transposing it, which may pay off over several subsequent computations
 - ▶ Finding the best layout for all multidimensional arrays is a NP-complete optimization problem [Mace, 1988]
 - **Example:** Recursive array layouts that preserve locality
 - ▶ Morton-order layout
 - ▶ Hierarchically tiled arrays
- In the best case, can make computations *cache-oblivious*
 - Performance largely independent of cache size

C. Kessler, IDA, Linköping universitet.

37

Loop Distribution (a.k.a. Loop Fission)



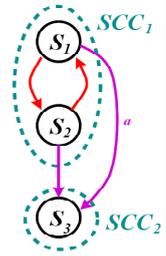
```

for (i=1; i<n; i++) {
S1:  a[i+1] = b[i-1] + c[i];
S2:  b[i]   = a[i] * k;
S3:  c[i]   = b[i] - 1;
}
    
```

↓ Loop distribution

```

for (i=1; i<n; i++) {
S1:  a[i+1] = b[i-1] + c[i];
S2:  b[i]   = a[i] * k;
}
for (i=1; i<n; i++)
S3:  c[i]   = b[i] - 1;
    
```



Safe if all statements forming a SCC in the dependence graph end up in the same loop.
 Forward (loop-carried) dep's are ok, but keep topological order.
 + often enables vectorization; better cache utilization of each loop.

C. Kessler, IDA, Linköping universitet.

38

Loop Fusion



```

for i from 1 to N do
  c[i] ← a[i] + b[i]
od
for j from 1 to N do
  d[j] ← a[j] * e[j]
od
    
```

fuse:

```

for i from 1 to N do
  c[i] ← a[i] + b[i]
  d[i] ← a[i] * e[i]
od
    
```

find second $a[i]$ in the cache or even in a register

For array a large enough, $a[i]$ will no longer be cached.

- Safe if neither loop carries a (backward) dependence.
- + locality: can convert inter-loop reuse to intra-loop reuse
- + larger basic blocks
- + reduce loop overhead

C. Kessler, IDA, Linköping universitet.

39

Loop Nest Flattening / Linearization



Flattens a multidimensional iteration space to a linear space:

```

for i from 0 to n-1 do
  for j from 0 to m-1 do
    iteration(i,j)
  od
od
    
```

linearize:

```

for k from 0 to m*n-1 do
  i ← k / m
  j ← k % m
  iteration(i,j)
od
    
```

- + larger iteration space, better for scheduling / load balancing
- overhead to reconstruct original iteration variables may be reduced by using *induction variables* i, j that are updated by accumulating additions instead of div and mod

C. Kessler, IDA, Linköping universitet.

40

Loop Unrolling



```

for i from 1 to 100 do
  a[i] ← a[i] + b[i]
od
    
```

unroll by 4:

```

for i from 1 to 100 step 4 do
  a[i]   ← a[i] + b[i]
  a[i+1] ← a[i+1] + b[i+1]
  a[i+2] ← a[i+2] + b[i+2]
  a[i+3] ← a[i+3] + b[i+3]
od
    
```

- + less overhead per useful operation
- + longer basic blocks for local optimizations (local CSE, local reg.-allocation, local scheduling, SW pipelining)
- longer code

C. Kessler, IDA, Linköping universitet.

41

Loop Unrolling with Unknown Upper Bound



```

i ← 1
while i+3 < N do
  a[i] ← a[i] + b[i]
  a[i+1] ← a[i+1] + b[i+1]
  a[i+2] ← a[i+2] + b[i+2]
  a[i+3] ← a[i+3] + b[i+3]
  i ← i + 4
od
while i < N do
  a[i] ← a[i] + b[i]
  i ← i + 1
od
    
```

used e.g. in BLAS

C. Kessler, IDA, Linköping universitet.

42

Loop Unroll-And-Jam

unroll the outer loop
and fuse the resulting inner loops:

```

for i from 1 to N do
  for j from 1 to N do
    a[i] ← a[i] + b[j]
  od
od
    
```

unroll&jam:

```

for i from 1 to N step 2 do
  for j from 1 to N do
    a[i] ← a[i] + b[j]
    a[i+1] ← a[i+1] + b[j]
  od
od
    
```

The same conditions as for loop interchange (for the two innermost loops after the unrolling step) must hold (for a formal treatment see [Allen/Kennedy'02, Ch. 8.4.1]).

- + increases reuse in inner loop
- + less overhead

C. Kessler, IDA, Linköpings universitet.

43

Loop Peeling

remove the first (or last) iteration of the loop
and clone the loop body for that iteration.

```

for i from 1 to N do
  a[i] ← (x+y)*b[i]
od
    
```

peel first iteration:

```

if N ≥ 1 then
  a[1] ← (x+y)*b[1]
  for i from 2 to N do
    a[i] ← (x+y)*b[i]
  od
fi
    
```

(Test on trip count can be removed if $N \geq 1$ is statically known.)

- + can enable loop fusion
- + may extract conditionals handling boundary cases from the loop
- longer code

C. Kessler, IDA, Linköpings universitet.

44

Index Set Splitting

Divide the *iteration space* into two portions.

```

for i from 1 to 100 do
  a[i] ← b[i] + c[i]
  if i > 10 then
    d[i] ← a[i] + a[i-10]
  fi
od
    
```

split after 10:

```

for i from 1 to 10 do
  a[i] ← b[i] + c[i]
od
for i from 11 to 100 do
  a[i] ← b[i] + c[i]
  d[i] ← a[i] + a[i-10]
od
    
```

- + removes condition evaluation in every iteration
- + factors out the parallelizable set of iterations
- longer code

C. Kessler, IDA, Linköpings universitet.

45

Loop Unswitching

```

for i from 1 to 100 do
  a[i] ← a[i] + b[i]
  if expression then
    d[i] ← 0
  fi
od
    
```

unswitch:

```

if expression then
  for i from 1 to 100 do
    a[i] ← a[i] + b[i]
    d[i] ← 0
  od
else
  for i from 1 to 100 do
    a[i] ← a[i] + b[i]
  od
fi
    
```

- + hoist loop-invariant control flow out of loop nest
- + no tests, no branches in loop body
 - larger basic blocks (see above), simpler software pipelining
- longer code

C. Kessler, IDA, Linköpings universitet.

46

Scalar Replacement

For (inner) loops accumulating a value in an array element
use a temporary scalar for the accumulator variable:

```

for i from 1 to N do
  for j from 1 to N do
    a[i] ← a[i] + b[j]
  od
od
    
```

scalar repl.:

```

for i from 1 to N step 2 do
  t ← a[i]
  for j from 1 to N do
    t ← t + b[j]
  od
  a[i] ← t
od
    
```

- + keep t in a register all the time
- + saves many costly memory accesses to $a[i]$

C. Kessler, IDA, Linköpings universitet.

47

Scalar Expansion / Array Privatization

promote a scalar temporary to an array to break a dependence cycle

```

if N ≥ 1
  allocate t'[1..N]
  for i from 1 to N do
    t ← a[i] + b[i]
    c[i] ← t + 1
  od
  t ← t' // if t live on exit
fi
    
```

expand scalar t:

- + removes the loop-carried antidependence due to t
 - can now parallelize the loop!

- needs more array space

Loop must be countable, scalar must not have upward exposed uses.

May also be done conceptually only, to enable parallelization:

just create one private copy of t for every processor = **array privatization**

C. Kessler, IDA, Linköpings universitet.

48

Idiom recognition and algorithm replacement

Traditional loop parallelization fails for loop-carried dep. with distance 1:

```
S0: s = 0;
    for (i=1; i<n; i++)
S1:   s = s + a[i];
S2: a[0] = c[0];
    for (i=1; i<n; i++)
S3:   a[i] = a[i-1] * b[i] + c[i];
```

↓ Idiom recognition (pattern matching)

```
S1': s = VSUM( a[1:n-1], 0 );
S3': a[0:n-1] = FOLR( b[1:n-1], c[0:n-1], mul, add );
```

↓ Algorithm replacement

```
S1'': s = par_sum( a, 0, n, 0 );
```

C. Kessler: Pattern-driven automatic parallelization. *Scientific Programming*, 1996.

A. Shafiee-Sarvestani, E. Hansson, C. Kessler: Extensible recognition of algorithmic patterns in DSP programs for automatic parallelization. *Int. J. on Parallel Programming*, 2013

C. Kessler, IDA, Linköpings universitet. 49



Concluding Remarks

Limits of Static Analyzability Outlook: Runtime Analysis and Parallelization

Christoph Kessler, IDA, Linköpings universitet, 2014.

Remark on static analyzability (1)

- Static dependence information is always a (safe) overapproximation of the real (run-time) dependences
 - Finding out the real ones exactly is statically undecidable!
 - If in doubt, a dependence must be assumed → may prevent some optimizations or parallelization
- One main reason for imprecision is **aliasing**, i.e. the program may have several ways to refer to the same memory location

Example: Pointer aliasing

```
void mergesort ( int* a, int n )
{
  ...
  mergesort ( a, n/2 );
  mergesort ( a + n/2, n-n/2 );
  ...
}
```

How could a static analysis tool (e.g., compiler) know that the two recursive calls read and write disjoint subarrays of a?

C. Kessler, IDA, Linköpings universitet. 51

Remark on static analyzability (2)

- Static dependence information is always a (safe) overapproximation of the real (run-time) dependences
 - Finding out the latter exactly is statically undecidable!
 - If in doubt, a dependence must be assumed → may prevent some optimizations or parallelization
- Another reason for imprecision are **statically unknown values** that imply whether a dependence exists or not

Example: Unknown dependence distance

```
// value of K statically unknown
for ( i=0; i<N; i++ )
{
  ...
  S: a[i] = a[i] + a[K];
  ...
}
```

Loop-carried dependence if $K < N$. Otherwise, the loop is parallelizable.

C. Kessler, IDA, Linköpings universitet. 52

Outlook: Runtime Parallelization

Sometimes parallelizability cannot be decided statically.

```
if is_parallelizable(...)
  forall i in [0..n-1] do // parallel version of the loop
    iteration(i);
  od
else
  for i from 0 to n-1 do // sequential version of the loop
    iteration(i);
  od
fi
```

The runtime dependence test `is_parallelizable(...)` itself may partially run in parallel.

C. Kessler, IDA, Linköpings universitet. 53



Run-Time Parallelization

Christoph Kessler, IDA, Linköpings universitet, 2014.

Goal of run-time parallelization



- Typical target: **irregular loops**

```
for ( i=0; i<n; i++)
    a[i] = f ( a[ g(i) ], a[ h(i) ], ... );
```

 - Array index expressions *g, h...* depend on run-time data
 - Iterations cannot be statically proved independent (and not either dependent with distance +1)
- **Principle:**
At runtime, inspect *g, h ...* to find out the real dependences and compute a schedule for partially parallel execution
 - Can also be combined with speculative parallelization

Overview



- **Run-time parallelization of irregular loops**
 - DOACROSS parallelization
 - Inspector-Executor Technique (shared memory)
 - Inspector-Executor Technique (message passing) *
 - Privatizing DOALL Test *
- **Speculative run-time parallelization of irregular loops ***
 - LRPD Test *
- **General Thread-Level Speculation**
 - Hardware support *

* = not covered in this course. See the references.

DOACROSS Parallelization



- Useful if loop-carried dependence distances are unknown, but often > 1
- Allow independent subsequent loop iterations to overlap
- Bilateral synchronization between really-dependent iterations

Example:

```
for ( i=0; i<n; i++)
    a[i] = f ( a[ g(i) ], ... );
```

↙

```
sh float aold[n];
sh flag done[n]; // flag (semaphore) array
forall i in 0..n-1 { // spawn n threads, one per iteration
    done[n] = 0;
    aold[i] = a[i]; // create a copy
}
forall i in 0..n-1 { // spawn n threads, one per iteration
    if (g(i) < i) wait until done[ g(i) ];
    a[i] = f ( a[ g(i) ], ... );
    set( done[i] );
}
else
    a[i] = f ( aold[ g(i) ], ... ); set done[i];
}
```

Inspector-Executor Technique (1)



- Compiler generates 2 pieces of customized code for such loops:

- **Inspector**
 - calculates values of index expression by simulating whole loop execution
 - ▶ typically, based on sequential version of the source loop (some computations could be left out)
 - computes implicitly the real iteration dependence graph
 - computes a parallel schedule as (greedy) wavefront traversal of the iteration dependence graph
 - ▶ all iterations in same wavefront are independent
 - ▶ schedule depth = #wavefronts = critical path length
- **Executor**
 - follows this schedule to execute the loop



Inspector-Executor Technique (2)



- **Source loop:**

```
for ( i=0; i<n; i++)
    a[i] = f ( a[ g(i) ], a[ h(i) ], ... );
```
- **Inspector:**

```
int wf[n]; // wavefront indices
int depth = 0;
for (i=0; i<n; i++)
    wf[i] = 0; // init.
for (i=0; i<n; i++) {
    wf[i] = max ( wf[ g(i) ], wf[ h(i) ], ... ) + 1;
    depth = max ( depth, wf[i] );
}
```



- Inspector considers only flow dependences (RAW), anti- and output dependences to be preserved by executor

Inspector-Executor Technique (3)



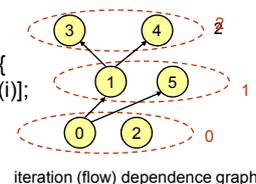
- **Example:**

```
for (i=0; i<n; i++)
    a[i] = ... a[ g(i) ] ...;
```

- **Executor:**

```
float aold[n]; // buffer array
aold[1:n] = a[1:n];
for (w=0; w<depth; w++)
    forall (i, 0, n, #) if (wf[i] == w) {
        a1 = (g(i) < i)? a[g(i)] : aold[g(i)];
        ... // similarly, a2 for h etc.
        a[i] = f ( a1, a2, ... );
    }
```

i	0	1	2	3	4	5
g(i)	2	0	2	1	1	0
wf[i]	0	1	0	2	2	1
g(i)<i?	no	yes	no	yes	yes	yes



Inspector-Executor Technique (4)



Problem: Inspector remains sequential – no speedup

Solution approaches:

- Re-use schedule over subsequent iterations of an outer loop if access pattern does not change
 - amortizes inspector overhead across repeated executions
- Parallelize the inspector using doacross parallelization [Saltz,Mirchandaney'91]
- Parallelize the inspector using sectioning [Leung/Zahorjan'91]
 - compute processor-local wavefronts in parallel, concatenate
 - trade-off schedule quality (depth) vs. inspector speed
 - Parallelize the inspector using bootstrapping [Leung/Z.'91]
 - Start with suboptimal schedule by sectioning, use this to execute the inspector → refined schedule

C. Kessler, IDA, Linköping universitet.

61

DF00100 Advanced Compiler Construction
TDDC86 Compiler optimizations and code generation



Thread-Level Speculation

Christoph Kessler, IDA, Linköping universitet, 2014.

Speculatively parallel execution

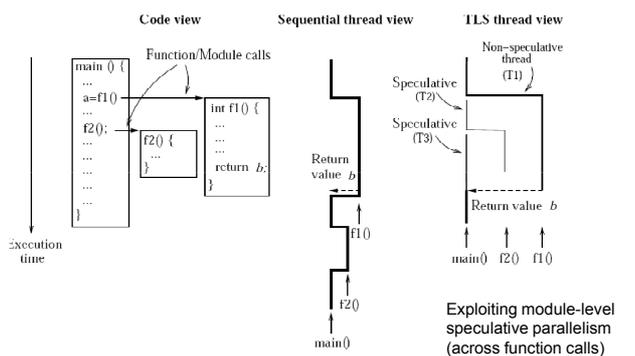


- For automatic parallelization of sequential code where dependences are hard to analyze statically
- Works on a **task graph**
 - constructed implicitly and dynamically
- **Speculate on:**
 - control flow, data independence, synchronization, values
 - We focus on thread-level speculation (TLS) for CMP/MT processors. Speculative instruction-level parallelism is not considered here.
- **Task:**
 - **statically:** Connected, single-entry subgraph of the control-flow graph
 - ▶ Basic blocks, loop bodies, loops, or entire functions
 - **dynamically:** Contiguous fragment of dynamic instruction stream within static task region, entered at static task entry

C. Kessler, IDA, Linköping universitet.

63

TLS Example

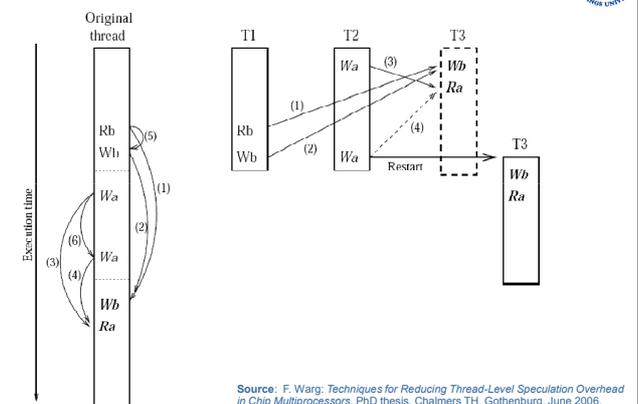


Source: F. Warg; *Techniques for Reducing Thread-Level Speculation Overhead in Chip Multiprocessors*. PhD thesis, Chalmers TH, Gothenburg, June 2006.

C. Kessler, IDA, Linköping universitet.

64

Data dependence problem in TLS



Source: F. Warg; *Techniques for Reducing Thread-Level Speculation Overhead in Chip Multiprocessors*. PhD thesis, Chalmers TH, Gothenburg, June 2006.

C. Kessler, IDA, Linköping universitet.

65

Speculatively parallel execution of tasks



- **Speculation on inter-task control flow**
 - After having assigned a task, predict its successor task and start it speculatively
- **Speculation on data independence**
 - For inter-task memory data (flow) dependences
 - ▶ conservatively: await write (memory synchronization, message)
 - ▶ speculatively: hope for independence and continue (execute the load)
- **Roll-back** of speculative results on mis-speculation (expensive)
 - When starting speculation, state must be buffered
 - Squash an offending task and all its successors, restart
- **Commit speculative results** when speculation resolved to correct
 - Task is retired

C. Kessler, IDA, Linköping universitet.

66

Selecting Tasks for Speculation



■ Small tasks:

- too much overhead (task startup, task retirement)
- low parallelism degree

■ Large tasks:

- higher misspeculation probability
- higher rollback cost
- many speculations ongoing in parallel may saturate the resources

■ Load balancing issues

- avoid large variation in task sizes

■ Traversal of the program's control flow graph (CFG)

- Heuristics for task size, control and data dep. speculation

C. Kessler, IDA, Linköpings universitet.

67

TLS Implementations



■ Software-only speculation

- for loops [Rauchwerger, Padua '94, '95]
- ...

■ Hardware-based speculation

- Typically, integrated in cache coherence protocols
- Used with multithreaded processors / chip multiprocessors for automatic parallelization of sequential legacy code
- If source code available, compiler may help e.g. with identifying suitable threads

C. Kessler, IDA, Linköpings universitet.

68

DF00100 Advanced Compiler Construction
TDDC86 Compiler optimizations and code generation



Questions?

Christoph Kessler, IDA,
Linköpings universitet, 2014.

Some references on Dependence Analysis Loop optimizations and Transformations



- H. Zima, B. Chapman: *Supercompilers for Parallel and Vector Computers*. Addison-Wesley / ACM press, 1990.
- M. Wolfe: *High-Performance Compilers for Parallel Computing* Addison-Wesley, 1996.
- R. Allen, K. Kennedy: *Optimizing Compilers for Modern Architectures*. Morgan Kaufmann, 2002.

Idiom recognition and algorithm replacement:

- C. Kessler: Pattern-driven automatic parallelization. *Scientific Programming* 5:251-274, 1996.
- A. Shafiee-Sarvestani, E. Hansson, C. Kessler: Extensible recognition of algorithmic patterns in DSP programs for automatic parallelization. *Int. J. on Parallel Programming*, 2013.

C. Kessler, IDA, Linköpings universitet.

70

Some references on run-time parallelization



- R. Cytron: Doacross: Beyond vectorization for multiprocessors. Proc. ICPP-1986
- D. Chen, J. Torrellas, P. Yew: An Efficient Algorithm for the Run-time Parallelization of DO-ACROSS Loops, Proc. IEEE Supercomputing Conf., Nov. 2004, IEEE CS Press, pp. 518-527
- R. Mirchandaney, J. Saltz, R. M. Smith, D. M. Nicol, K. Crowley: Principles of run-time support for parallel processors, Proc. ACM Int. Conf. on Supercomputing, July 1988, pp. 140-152.
- J. Saltz and K. Crowley and R. Mirchandaney and H. Berryman: Runtime Scheduling and Execution of Loops on Message Passing Machines, *Journal on Parallel and Distr. Computing* 8 (1990): 303-312.
- J. Saltz, R. Mirchandaney: The preprocessed doacross loop. Proc. ICPP-1991 Int. Conf. on Parallel Processing.
- S. Leung, J. Zahorjan: Improving the performance of run-time parallelization. Proc. ACM PPoPP-1993, pp. 83-91.
- Lawrence Rauchwerger, David Padua: The Privatizing DOALL Test: A Run-Time Technique for DOALL Loop Identification and Array Privatization. Proc. ACM Int. Conf. on Supercomputing, July 1994, pp. 33-45.
- Lawrence Rauchwerger, David Padua: The LRPD Test: Speculative Run-Time Parallelization of Loops with Privatization and Reduction Parallelization. Proc. ACM SIGPLAN PLDI-95, 1995, pp. 218-232.

C. Kessler, IDA, Linköpings universitet.

71

Some references on speculative execution / parallelization



- T. Vijaykumar, G. Sohi: Task Selection for a Multiscalar Processor. Proc. MICRO-31, Dec. 1998.
- J. Martinez, J. Torrellas: Speculative Locks for Concurrent Execution of Critical Sections in Shared-Memory Multiprocessors. Proc. WMPJ at ISCA, 2001.
- F. Warg and P. Stenström: Limits on speculative module-level parallelism in imperative and object-oriented programs on CMP platforms. Pr. IEEE PACT 2001.
- P. Marcuello and A. Gonzalez: Thread-spawning schemes for speculative multithreading. Proc. HPCA-8, 2002.
- J. Steffan et al.: Improving value communication for thread-level speculation. HPCA-8, 2002.
- M. Cintra, J. Torrellas: Eliminating squashes through learning cross-thread violations in speculative parallelization for multiprocessors. HPCA-8, 2002.
- Fredrik Warg and Per Stenström: Improving speculative thread-level parallelism through module run-length prediction. Proc. IPDPS 2003.
- F. Warg: *Techniques for Reducing Thread-Level Speculation Overhead in Chip Multiprocessors*. PhD thesis, Chalmers TH, Gothenburg, June 2006.
- T. Ohsawa et al.: Pinot: Speculative multi-threading processor architecture exploiting parallelism over a wide range of granularities. Proc. MICRO-38, 2005.

C. Kessler, IDA, Linköpings universitet.

72