# Automatic Translation of Fortran Programs to Vector Form

Randy Allen and Ken Kennedy

# The problem

- New (as of 1987) vector machines such as the Cray-1 have proven successful

- Most Fortran code is written sequentially, using loops

- Can we exploit parallelism without rewriting everything?

# Compiler Vectorization

- Idea: Compiler detects parallelism and automatically converts loops

- No need to rewrite code or learn new language

- Opportunities for parallelism are subtle and difficult to detect

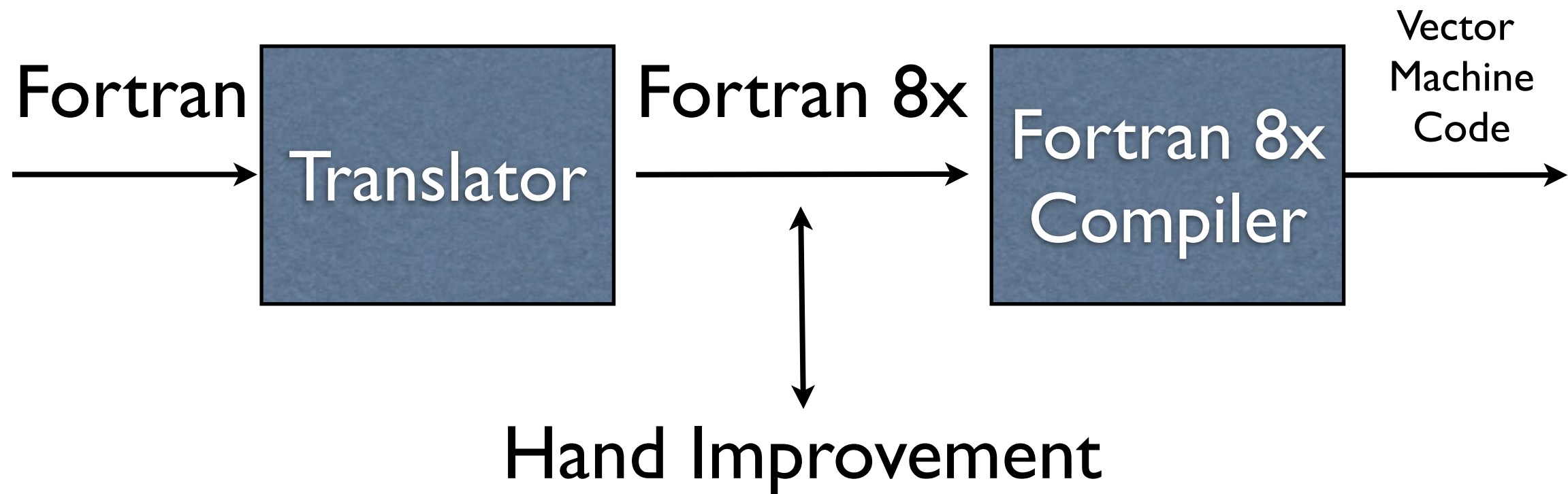- Programmers need to tweak code into forms the optimizer can recognize

# Explicit Vector Instructions

- Idea: Add forms to Fortran 8x to specify parallel operations

- Avoid writing sequential code in the first place

- Programmers better understand what parallelism opportunities exist

- Code must be rewritten

# Source-to-Source Translation

- Idea: Automatically convert existing sequential Fortran into parallel Fortran 8x

- Translation only occurs once, so more expensive transformations are practical

- Programmers can add any needed parallelism the translator misses

# Parallel Fortran Converter

# Vector Operations in Fortran 8x

- Note: As of this paper, Fortran 8x is still theoretical

- Vectors and arrays may be treated as aggregates: $X = Y + Z$

- Arithmetic operators are applied point-wise

- Scalars are treated as same-valued vectors

- All arrays must have the same dimensions

# Simultaneity

- Array assignment (e.g. X = Y) is *simultaneous*. All of Y is fetched before X is stored

- X = X / X(2) uses the value for X(2) prior to the assignment, even though X(2) will be assigned to

- Equivalent to using a temporary variable

# Array Sections

- *Triplet* notation allows reference to parts of arrays

- `A(1:100, I) = B(J, 1:100)` assigns 100 elements from row J of B to column I of A

- Third element of triple specifies *stride*: A(2:100:2) references first 50 even subscript positions

# Array Identification

- Specifies a named mapping to an array

- `IDENTIFY /1:M/ D(I) = C(I, I + 1)` defines D as the superdiagonal of C

- D is just an alias; it has no storage

# Conditional Assignment

- `WHERE(A .GT. 0.0) A = A + B` indicates that only positive elements of A will be modified

- Errors in evaluating the right-hand side must be ignored when the predicate fails

- E.g., `WHERE(A .NE. 0.0) B = B/A`

# Library Functions

- Mathematical functions (SIN, SQRT, etc.) are extended to operate on arrays

- New intrinsic array operations: DOTPRODUCT, TRANSPOSE

- `SEQ(1,N)` returns an index array

- Reductions operations, e.g. SUM

# Translation Process

```
          DO 20 I = 1, 100
$S_1$         KI = I
          DO 10 J = 1, 300, 3
$S_2$          KI = KI + 2
$S_3$          U(J) = U(J) * W(KI)
$S_4$          V(J + 3) = V(J) + W(KI)
     10    CONTINUE
     20  CONTINUE
```

- Goal: transform $S_3$ and $S_4$ into vector instructions and remove them from the inner loop

- Only possible if there is no semantic difference

# Simple Case

```
DO 10 I = 1, 100
     X(I) = X(I) + Y(I)
10   CONTINUE
```

- Easily becomes `X(1:100) = X(1:100) + Y(1:100)`

```
DO 10 I = 1, 100
     X(I + 1) = X(I) + Y(I)
10   CONTINUE
```

- Cannot be converted, because each iteration depends on the previous

  - Known as a *recurrence*

# Dependency Detection

- To distinguish parallel and non-parallel loops, translator must detect self-dependent statements

- First, code is *normalized* to make this test feasible

# DO-Loop Normalization

```
        DO 20 I = 1, 100
           KI = I
           DO 10 j = 1, 100
              KI = KI + 2
              U(3 * j − 2) = U(3 * j − 2) * W(KI)
              V(3 * j + 1) = V(3 * j − 2) + W(KI)
     10      CONTINUE
S₆        J = 301
     20   CONTINUE
```

- Convert induction variables to iterate from 1 by steps of 1

- Here, J has been replaced by j

- $S_6$ added to preserve post-condition

# Induction Variable Substitution

```
          DO 20 I = 1, 100
              KI = I
              DO 10 j = 1, 100
                  U(3 * j − 2) = U(3 * j − 2) * W(KI + 2 * j)
                  V(3 * j + 1) = V(3 * j − 2) + W(KI + 2 * j)
     10       CONTINUE
              KI = KI + 200
              J = 301
     20   CONTINUE
```

- Convert all subscripts to linear functions of induction variables

- KI has been removed from loop and replaced by its initial value plus its increments

- KI updated post-loop with final value

- Note: repeated addition replaced by multiplication

# Dead Statement Elimination

```
        DO 20 I = 1, 100
          DO 10 j = 1, 100
```
$S_3$ $\qquad$ $U(3 * j - 2) = U(3 * j - 2) * W(I + 2 * j)$
$S_4$ $\qquad$ $V(3 * j + 1) = V(3 * j - 2) + W(I + 2 * j)$
```
      10    CONTINUE
      20  CONTINUE
```
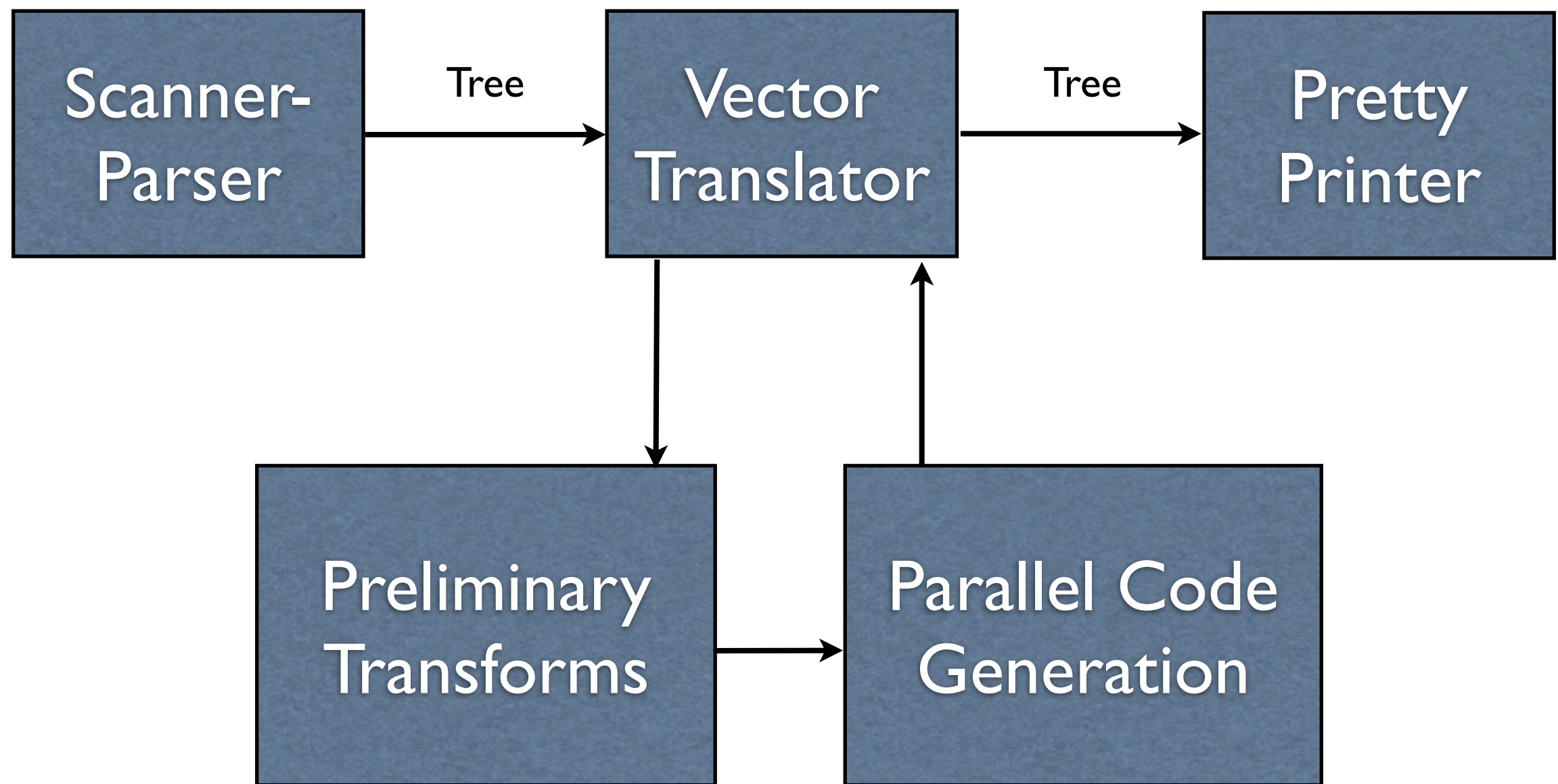
- Assuming J and KI aren't used outside the loop, their final values can be discarded

- Since they also aren't used within the loop, they can be removed entirely

# Vector Code Generation

```
        DO 20 I = 1, 100
S₃          U(1:298:3) = U(1:298:3) * W(I − 2:I + 200:2)
            DO 10 j = 1, 100
S₄              V(3 * j + 1) = V(3 * j − 2) + W(I + 2 * j)
    10      CONTINUE
    20   CONTINUE
```

- Dependency analysis shows $S_4$ depends on itself, but $S_3$ does not

- Therefore, $S_3$ can be vectorized and moved out of the loop

# Translation Process

# Dependence Analysis

- $S_2$ depends on $S_1$ if some execution of $S_2$ uses a value from a previous $S_1$

- Self-dependence can only arise in loops

# Dependency in Loops

```
      DO 10 J = 1, N
         X(J) = X(J) + C
 10   CONTINUE
```

No dependency

```
      DO 10 J = 1, N − 1
         X(J + 1) = X(J) + C
 10   CONTINUE
```

Recurrence

# Dependency in Loops

$$\text{(*)} \quad \begin{array}{l} \text{DO } 10 \text{ i} = 1, \text{ N} \\ \quad \text{X}(f(i)), = \textbf{F}(\text{X}(g(i))) \\ 10 \quad \text{CONTINUE} \end{array} \quad \text{General Form}$$

- (*) depends on itself iff there exist $i_1$, $i_2$ such that $1 \leq i_1 < i_2 \leq N$ and $f(i_1) = g(i_2)$

- Most often, $f$ and $g$ are linear in $i$

- $ax + by = n$ has a linear solution iff $gcd(a,b) \mid n$

- $f(i) = a_0 + a_1 i$; $g(i) = b_0 + b_1 i$

- (*) depends on itself only if $gcd(a_1, b_1) \mid b_0 - a_0$

# Dependency in Loops

- Unfortunately, $gcd(a_1, b_1)$ is commonly 1

- More sophisticated techniques are needed

COROLLARY 3 (BANERJEE INEQUALITY). *If $f(x) = a_0 + a_1 x$ and $g(y) = b_0 + b_1 y$ then statement (\*) depends on itself only if*

$$-b_1 - (a_1^- + b_1)^+(N - 2) \leq b_0 + b_1 - a_0 - a_1 \leq -b_1 + (a_1^+ - b_1)^+(N - 2).$$

- Even these only provide *necessary* conditions for dependence

- Multiple loops are harder still

# Indirect Dependence

$$
\begin{array}{ll}
& \text{DO 10 I = 1, 100} \\
S_1 & \quad \text{T(I) = A(I) * B(I)} \\
S_2 & \quad \text{S(I) = S(I) + T(I)} \\
S_3 & \quad \text{A(I + 1) = S(I) + C(I)} \\
10 & \text{CONTINUE}
\end{array}
$$

- $S_1$, $S_2$, and $S_3$ all depend indirectly on themselves

# Types of Dependency

- We say $S_2$ depends on $S_1$ if one of these conditions hold

- *True dependence*: $S_2$ uses the output of $S_1$

- *Antidependence*: $S_1$ would use the output of $S_2$ if they were reversed

- *Output dependence*: $S_2$ recalculates the output of $S_1$

# Loop-Related Dependency

- *Loop carried dependence*: one statement stores to a location; another statement reads from that location in a *later* iteration

- *Loop independent dependence*: one statement stores to a location; another statement reads from that location in the *same* iteration

- Self-dependence is a special case of loop carried dependence

# Testing Procedure

- Test each pair of statements for dependence, building a dependence relation $D$

- Compute the transitive closure $D^+$

- Execute statements which do not depend on themselves in $D^+$ in parallel

- Execution order must be consistent with $D^+$

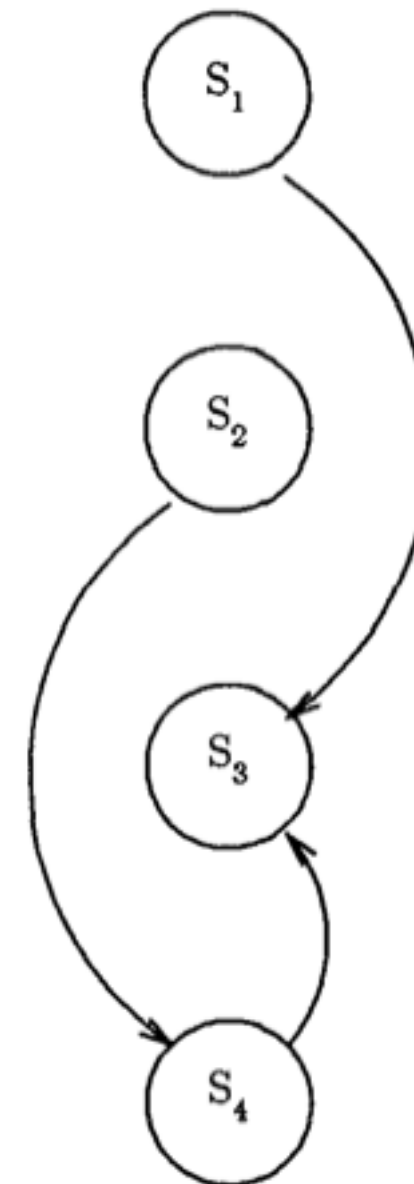- Reduce cycles to $\pi$-blocks; the resulting graph is acyclic

# Example

## This…

```
        DO 10 I = 1, 99
S₁          X(I) = I
S₂          B(I) = 100-I
   10   CONTINUE
        DO 20 I = 1, 99
S₃          A(I) = F(X(I))
S₄          X(I + 1) = G(B(I))
   20   CONTINUE
```

$S_1$
$S_2$

$X(I) = I$

$X(I+1) = G(B(I))$

## Becomes…

$$X(1:99) = SEQ(1, 99, 1)$$
$$B(1:99) = SEQ(99, 1, -1)$$
$$X(2:100) = G(B(1:99))$$
$$A(1:99) = F(X(1:99))$$

## Note: $S_4$ precedes $S_3$

$S_1$

$S_2$

$S_3$

$S_4$

# Multiple Loops

- Important to note which loop carries the dependence

- We can define a maximum depth where a given dependence occurs

- Loop independent dependencies have infinite depth

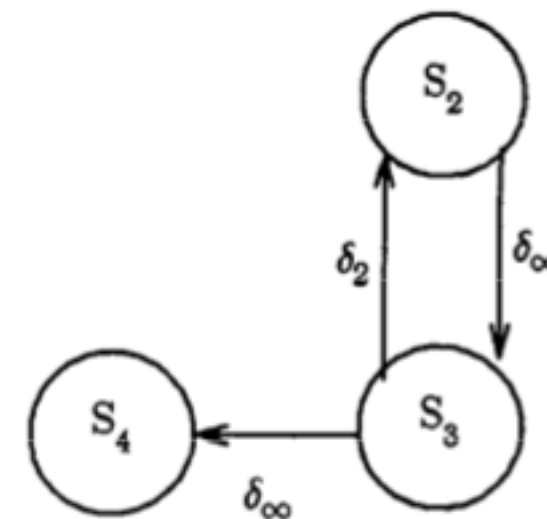- Dependency arcs are labeled with depth and type

# Example

# Example

```
        DO 30 I = 1, 100
           DO 20 J = 1, 100

              code for S₂, S₃
              generated at lower levels

   20      CONTINUE
S₄         Y(I + 1:I + 100) = A(2:101, N)
   30   CONTINUE
S₁      X(1:100) = Y(1:100) + 10
```

```
        DO 30 I = 1, 100
           DO 20 J = 1, 100
S₂            B(J) = A(J, N)
S₃            A(J + 1, 1:100) = B(J) + C(J, 1:100)
   20      CONTINUE
S₄         Y(I + 1:I + 100) = A(2:101, N)
   30   CONTINUE
S₁      X(1:100) = Y(1:100) + 10
```

# Further Techniques

- *Loop interchange*: move recurrences to outer loops

- *Recurrence breaking*: antidependent and output dependent single-statement recurrences can be ignored

- *Thresholds*: recurrences may permit partial vectorization

# Conditional Statements

Initial code

```
      DO 100 I = 1, N
         IF (A(I) .LE. 0) GOTO 100
         A(I) = B(I) + 3
100   CONTINUE
```

Convert to data dependency

```
      DO 100 I = 1, N
         BR1(I) = A(I) .LE. 0
         IF (.NOT. BR1(I)) A(I) = B(I) + 3
100   CONTINUE
```

Vectorize

```
BR1(1:N) = A(1:N) .LE. 0
WHERE (.NOT. BR1(1:N)) A(1:N) = B(1:N) + 3
```
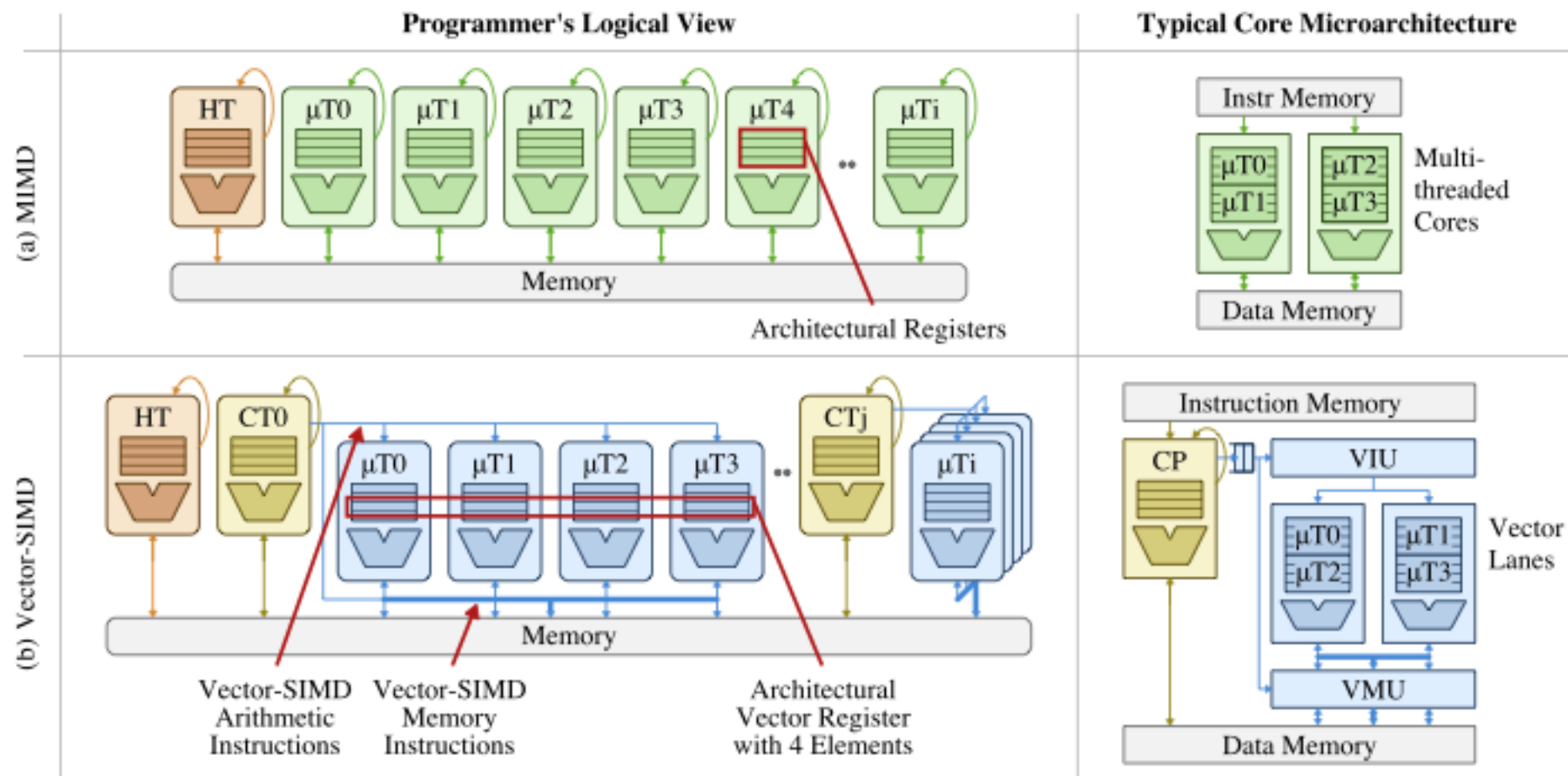
# Implementation

- Initial work based on PARAFRASE

- PFC is ~25,000 lines of PL/I

- Implements most of the translations discussed in the paper

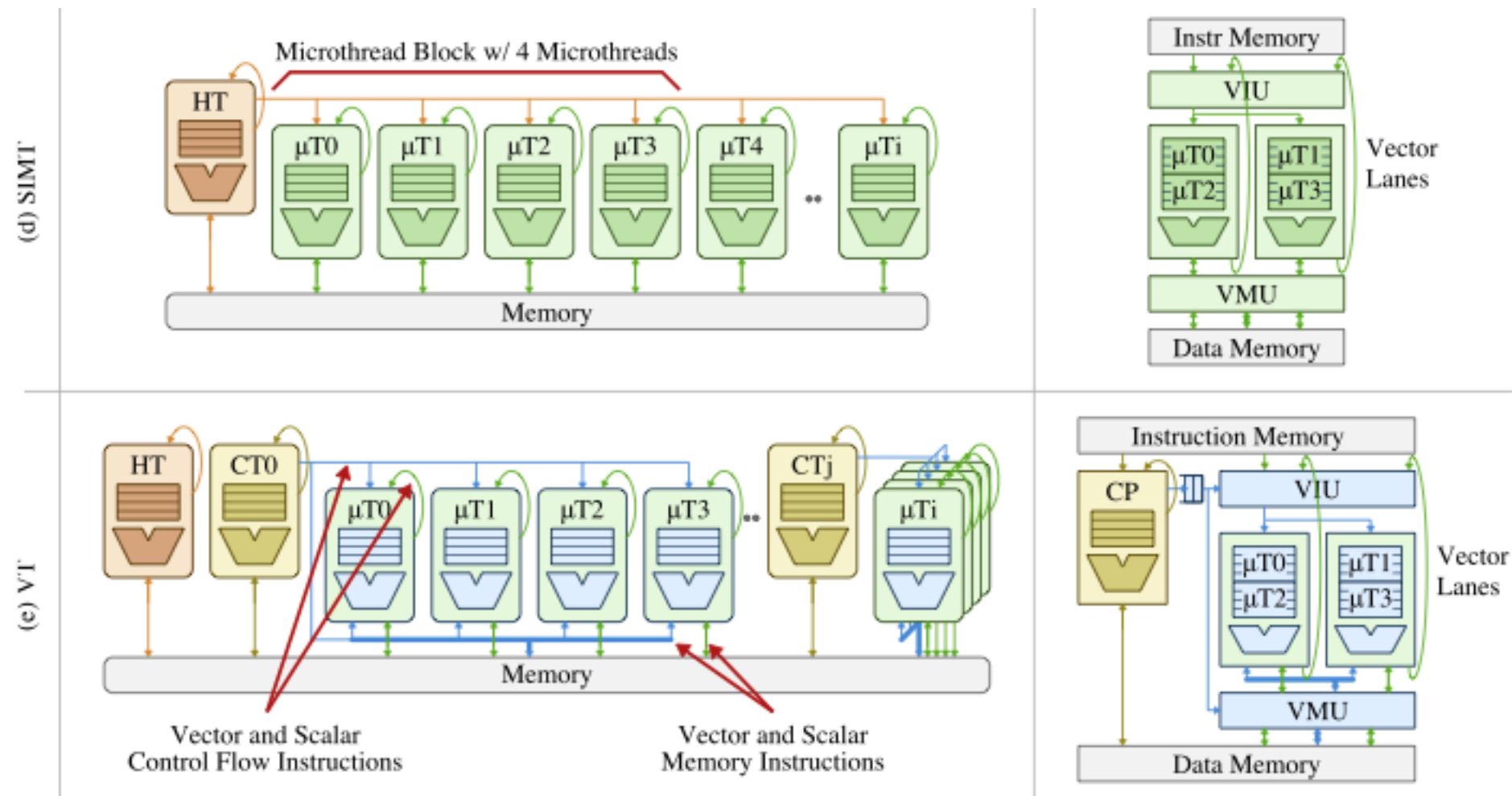- Runs their test case in 1 min on a 3 MB machine

# Exploring the tradeoffs between programmability and efficiency in data-parallel accelerators

Yunsup Lee, Rimas Avizienis, Alex Bishara, Richard Xia,
Derek Lockhart, Christopher Batten and Krste Asanović

# MIMD vs SIMD

# Two Hybrid Approaches

# SIMT

- Combines MIMD's logical view with vector-SIMD's microarchitecture

- VIU executes multiple µTs using SIMD as long as they proceed on the same control path

- VIU uses masks to selectively disable inactive µTs on different paths

# VT

- HT manages CTs; CTs manage µTs

- Vector-fetch instruction indicates scalar instructions to be executed by µTs

- VIU operates µTs in SIMD manner, but scalar branch can cause divergence

# Irregular Control Flow Example

```
for ( i = 0; i < n; i++ )
  if ( A[i] > 0 )
    C[i] = x * A[i] + B[i];
```

```
1   load      x, x_ptr
2 loop:
3   setvl     vlen, n
4   load.v    VA, a_ptr
5   load.v    VB, b_ptr
6   cmp.gt.v  VF, VA, 0
7   mul.sv    VT, x, VA,  VF
8   add.vv    VC, VT, VB, VF
9   store.v   VC, c_ptr,  VF
10  add       a_ptr, vlen
11  add       b_ptr, vlen
12  add       c_ptr, vlen
13  sub       n, vlen
14  br.neq    n, 0, loop
```
(b) Vector-SIMD

```
1   br.gte tidx, n, done
2   add    a_ptr, tidx
3   load   a, a_ptr
4   br.eq  a, 0, done
5   add    b_ptr, tidx
6   add    c_ptr, tidx
7   load   x, x_ptr
8   load   b, b_ptr
9   mul    t, x, a
10  add    c, t, b
11  store  c, c_ptr
12 done:
```
(c) SIMT

```
1   load      x, x_ptr
2   mov.sv    VZ, x
3 loop:
4   setvl     vlen, n
5   load.v    VA, a_ptr
6   load.v    VB, b_ptr
7   mov.sv    VD, c_ptr
8   fetch.v   ut_code
9   add       a_ptr, vlen
10  add       b_ptr, vlen
11  add       c_ptr, vlen
12  sub       n, vlen
13  br.neq    n, 0, loop
14  ...
15 ut_code:
16  br.eq     a, 0, done
17  mul       t, z, a
18  add       c, t, b
19  add       d, tidx
20  store     c, d
21 done:
22  stop
```
(d) VT

# Summary

- Vector-based microarchitectures more area and energy efficient than scalar-based

- Maven (VT) more efficient and easier to program than vector-SIMD

- Suggestion that VT more efficient but harder to program than SIMT