

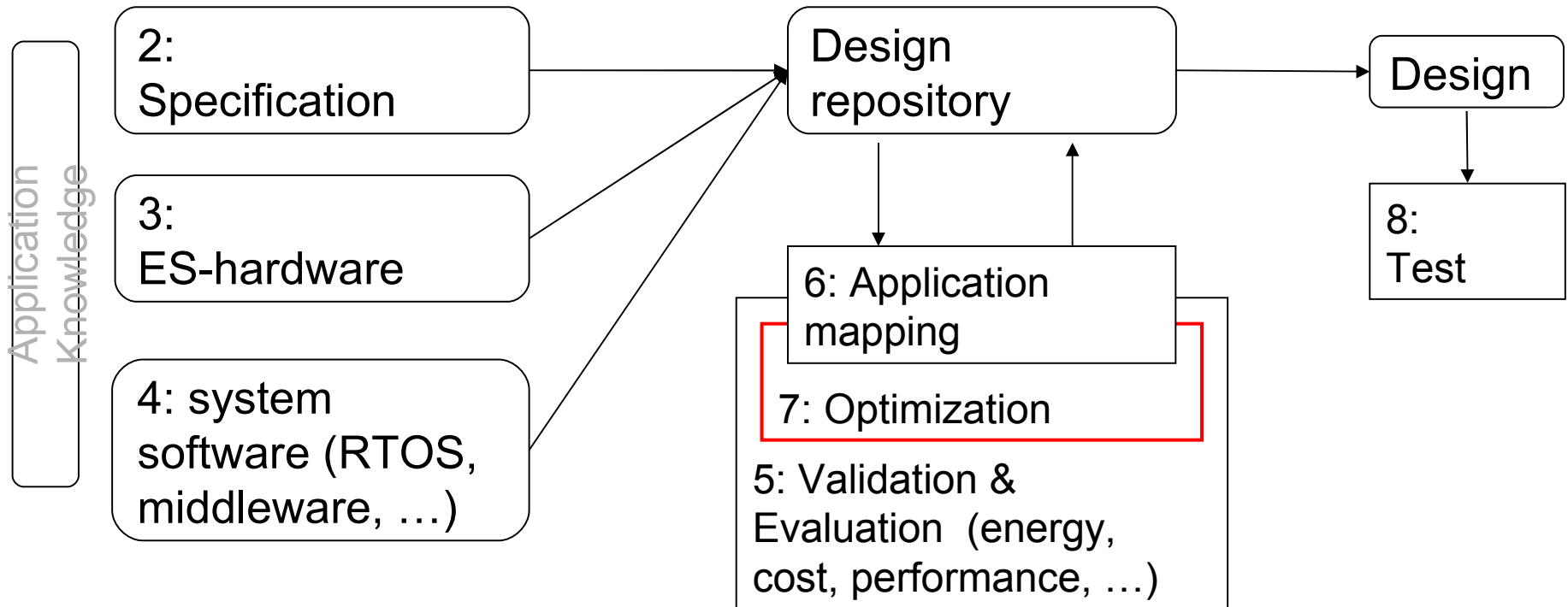
Optimizations

Peter Marwedel
TU Dortmund
Informatik 12
Germany

2010/01/13



Structure of this course



Numbers denote sequence of chapters

Improving predictability for caches

- Loop caches
- Mapping code to less used part(s) of the index space
- Cache locking/freezing
- Changing the memory allocation for code or data
- Mapping pieces of software to specific ways

Methods:

- Generating appropriate way in software
- Allocation of certain parts of the address space to a specific way
- Including way-identifiers in virtual to real-address translation

 “Caches behave almost like a scratch pad”

Code Layout Transformations (1)

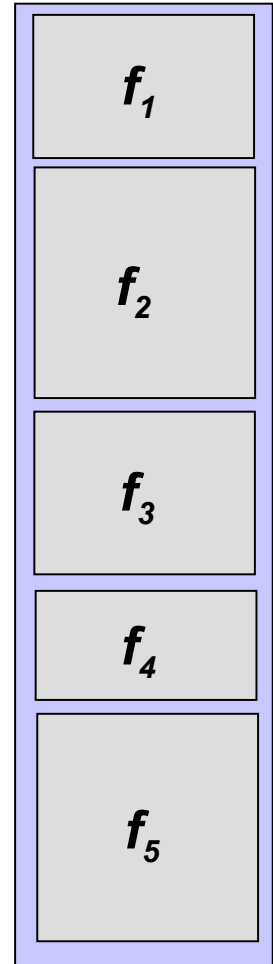
Execution counts based approach:

- Sort the functions according to execution counts (1100)
 $f_4 > f_1 > f_2 > f_5 > f_3$
- Place functions in decreasing order of execution counts (900)

(400)

(2000)

(700)



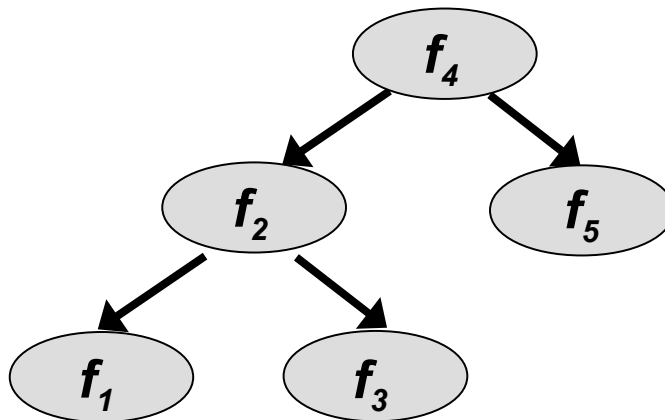
[S. McFarling: Program optimization for instruction caches, 3rd International Conference on Architecture Support for Programming Languages and Operating Systems (ASPLOS), 1989]

Code Layout Transformations (2)

Execution counts based approach:

- Sort the functions according to execution counts
 $f_4 > f_1 > f_2 > f_5 > f_3$
- Place functions in decreasing order of execution counts

Transformation increases spatial locality.
Does not take in account calling order



(2000)

f_4

(1100)

f_1

(900)

f_2

(700)

f_5

(400)

f_3

Code Layout Transformations (3)

Call-Graph Based Algorithm:

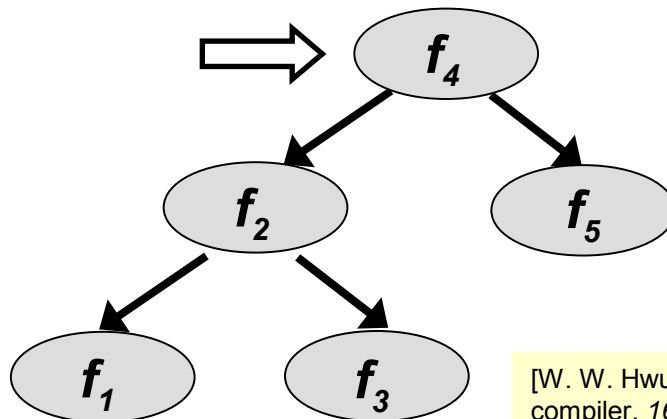
- Create weighted call-graph.
- Place functions according to weighted depth-first traversal.

$$f_4 > f_2 > f_1 > f_3 > f_5$$

Increases spatial locality.

(2000)

f_4



[W. W. Hwu et al.: Achieving high instruction cache performance with an optimizing compiler, 16th Annual International Symposium on Computer Architecture, 1989]

Code Layout Transformations (3)

Call-Graph Based Algorithm:

- Create weighted call-graph.
- Place functions according to weighted depth-first traversal.

$$f_4 > f_2 > f_1 > f_3 > f_5$$

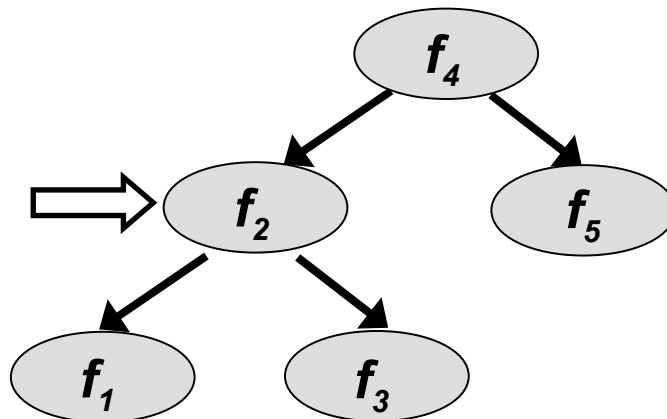
Increases spatial locality.

(2000)

(900)

f_4

f_2



Code Layout Transformations (4)

Call-Graph Based Algorithm:

- Create weighted call-graph.
- Place functions according to weighted depth-first traversal.

$$f_4 > f_2 > f_1 > f_3 > f_5$$

Increases spatial locality.

(2000)

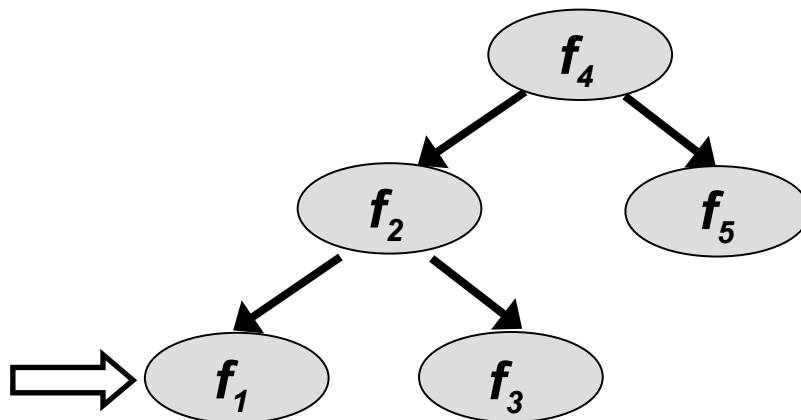
f_4

(900)

f_2

(1100)

f_1



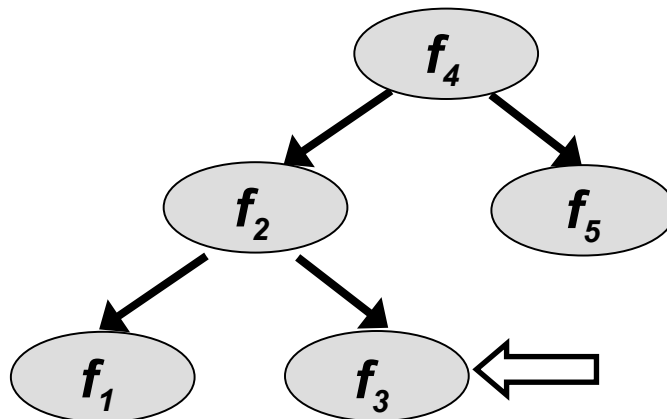
Code Layout Transformations (5)

Call-Graph Based Algorithm:

- Create weighted call-graph.
- Place functions according to weighted depth-first traversal.

$$f_4 > f_2 > f_1 > f_3 > f_5$$

Increases spatial locality.



(2000)

f_4

(900)

f_2

(1100)

f_1

(400)

f_3

Code Layout Transformations (6)

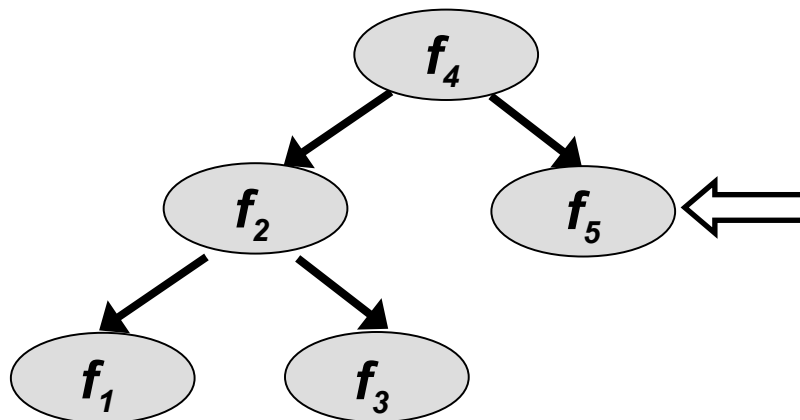
Call-Graph Based Algorithm:

- Create weighted call-graph.
- Place functions according to weighted depth-first traversal.

$$f_4 > f_2 > f_1 > f_3 > f_5$$

- Combined with placing frequently executed traces at the top of the code space for functions.

Increases spatial locality.



(2000)

f_4

(900)

f_2

(1100)

f_1

(400)

f_3

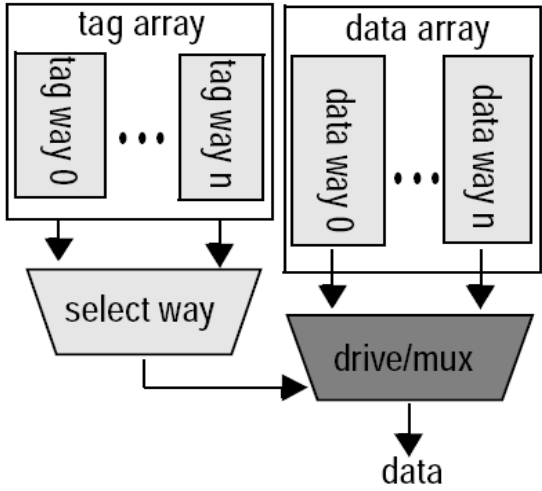
(700)

f_5

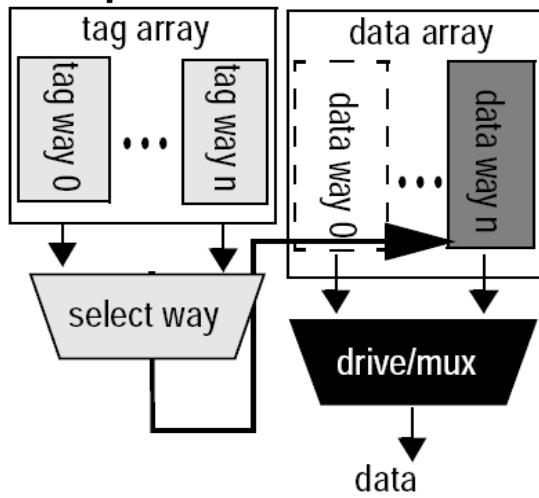
Way prediction/selective direct mapping

Timing order: 1st step 2nd step 3rd step No activity

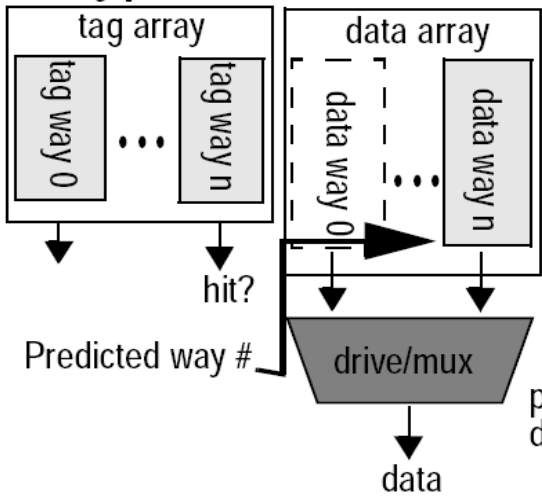
a: Conventional parallel access



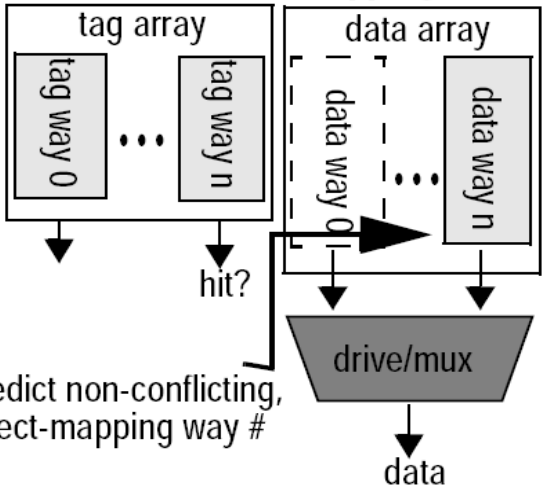
b: Sequential access



c: Way-prediction



d: Selective direct-mapping



[M. D. Powell, A. Agarwal, T. N. Vijaykumar, B. Falsafi, K. Roy: Reducing Set-Associative Cache Energy via Way-Prediction and Selective Direct-Mapping, *MICRO-34*, 2001]

Hardware organization for way prediction

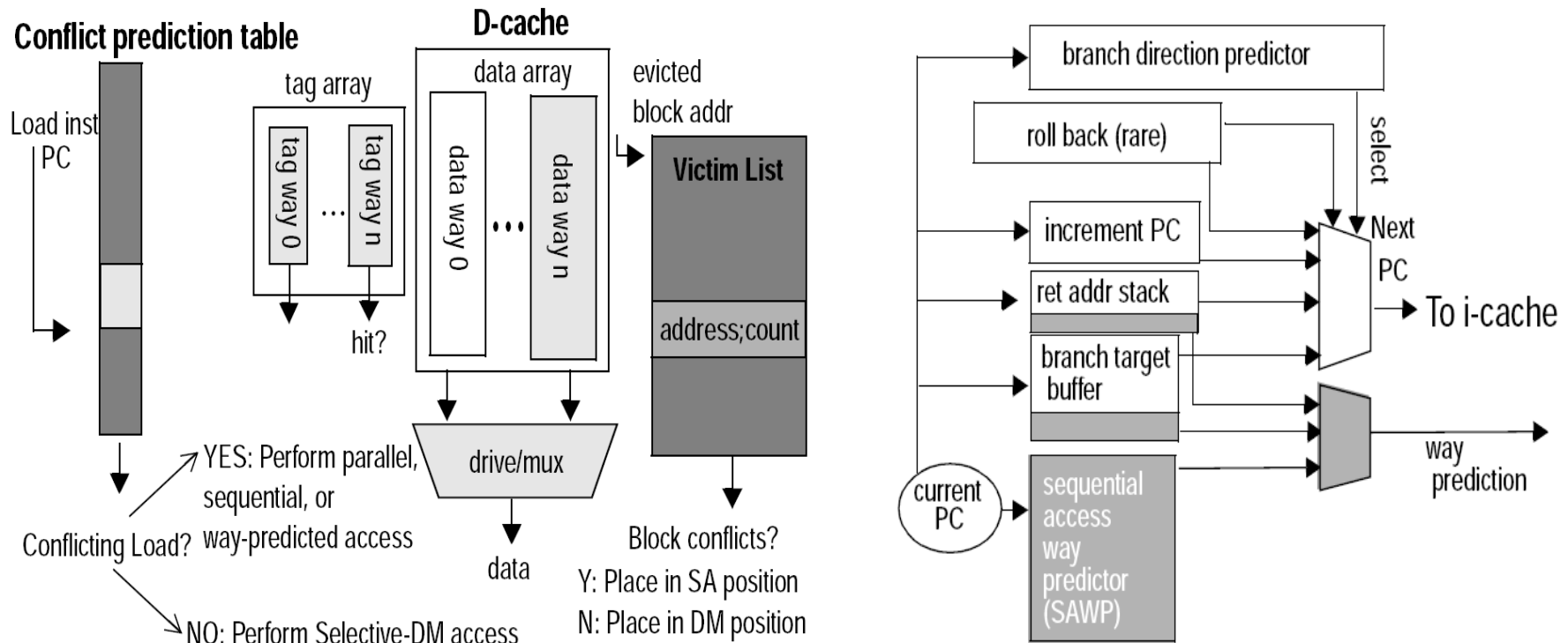


FIGURE 3: Fetch and i-cache access mechanism.

Results for the paper on way prediction (1)

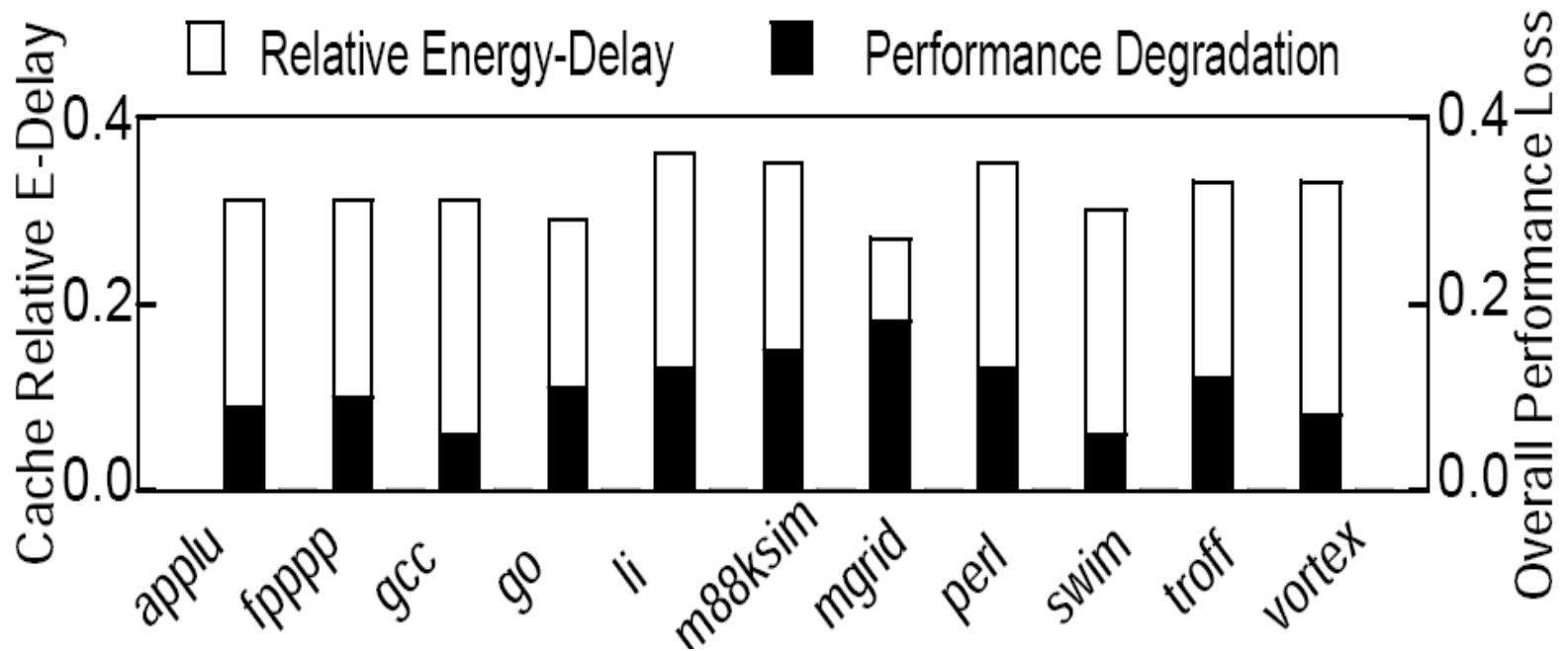
System configuration parameters

| | |
|--------------------------------------|-----------------------------------|
| Instruction issue & decode bandwidth | 8 issues per cycle |
| L1 I-Cache | 16K, 4-way, 1 cycle |
| Base L1 D-Cache | 16K, 4-way, 1 or 2 cycles, 2ports |
| L2 cache | 1M, 8-way, 12 cycle latency |
| Memory access latency | 80 cycles+4 cycles per 8 bytes |
| Reorder buffer size | 64 |
| LSQ size | 32 |
| Branch predictor | 2-level hybrid |

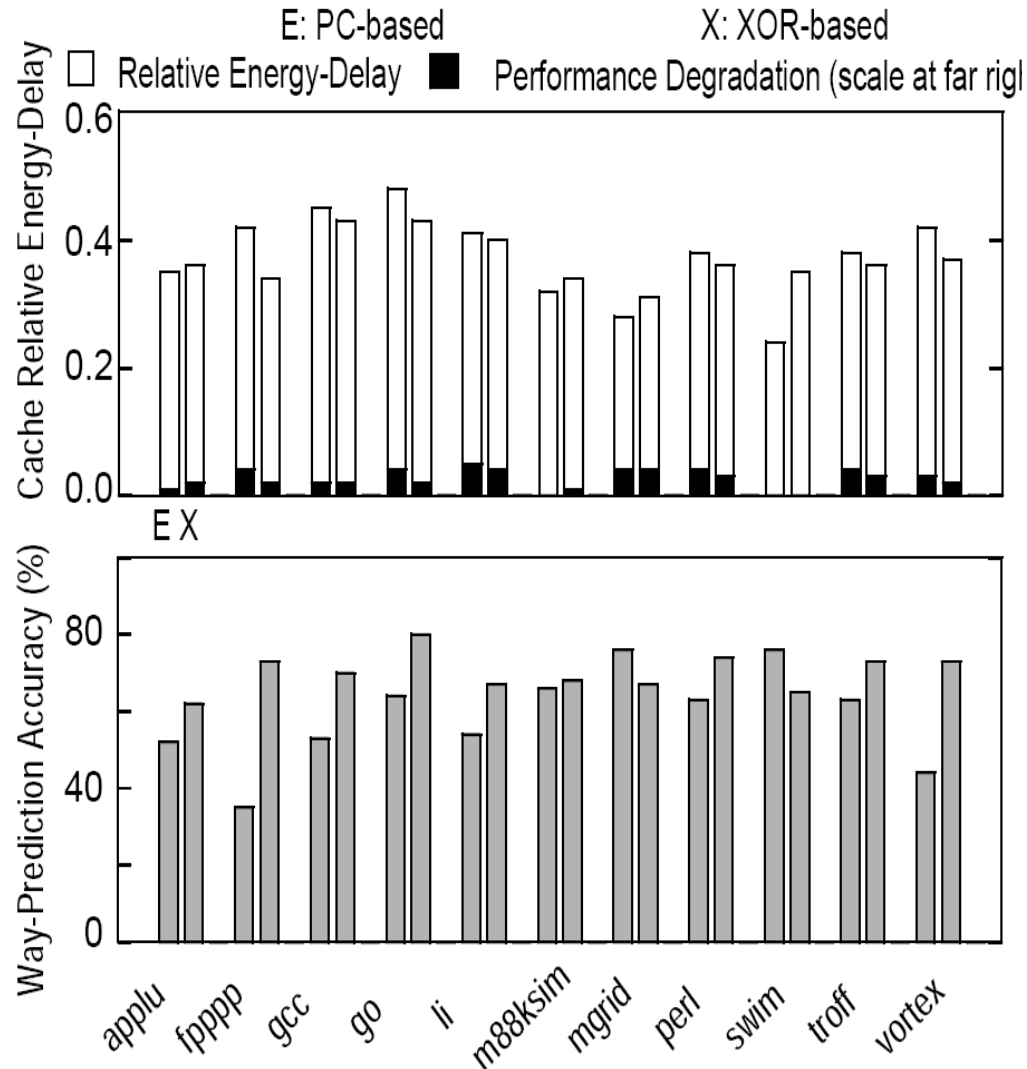
Cache energy and prediction overhead

| Energy component | Relative energy |
|---|-----------------|
| Parallel access cache read (4 ways read) | 1.00 |
| 1 way read | 0.21 |
| Cache write | 0.24 |
| Tag array energy (incl. in the above numbers) | 0.06 |
| 1024x4bit prediction table read/write | 0.007 |

Results for the paper on way prediction (2)



Results for the paper on way prediction (2)



The offset assignment problem

Peter Marwedel
TU Dortmund
Informatik 12
Germany

2010/01/13

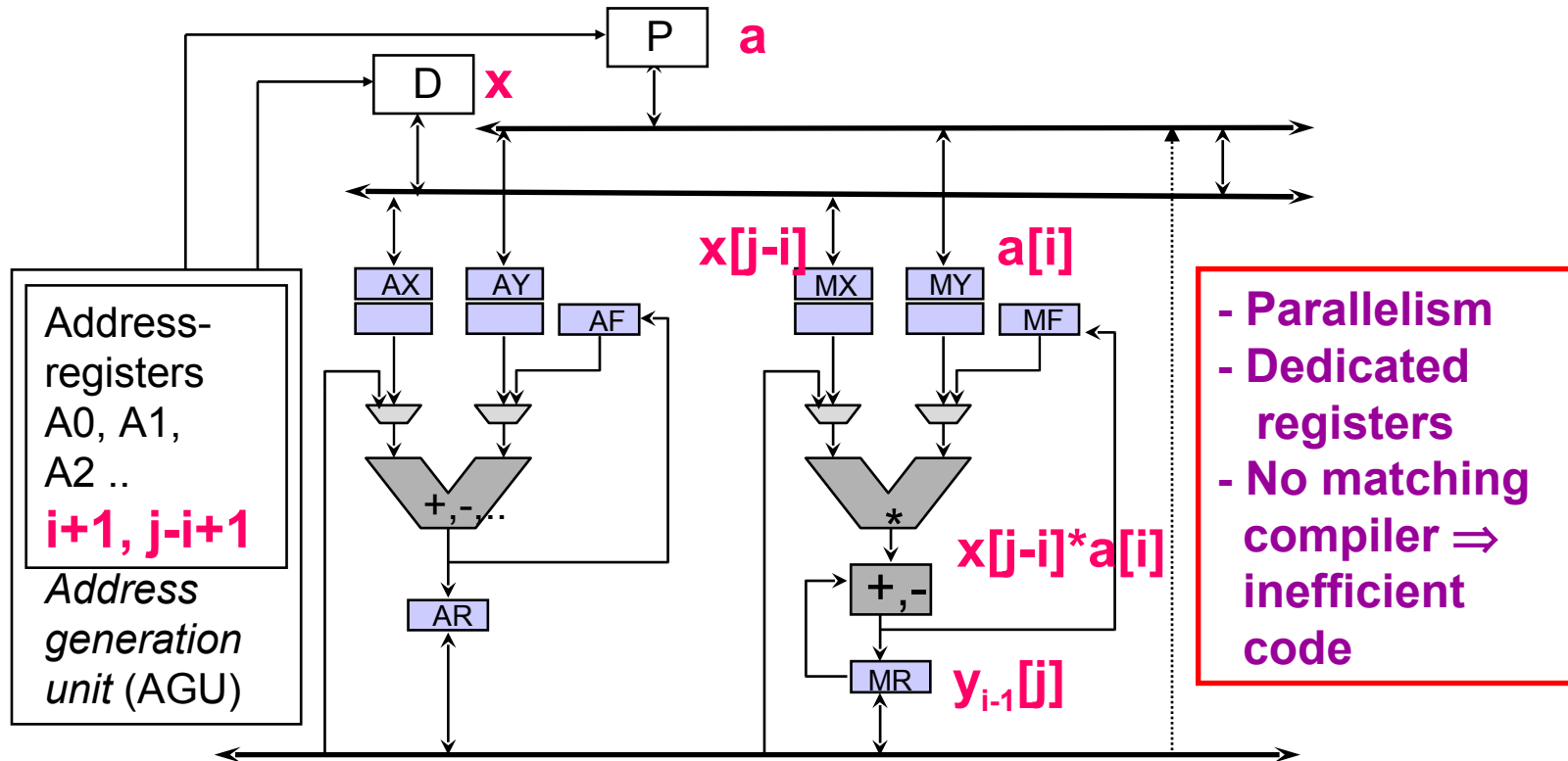


Reason for compiler-problems: Application-oriented Architectures

Application: u.a.: $y[j] = \sum_{i=0}^n x[j-i] * a[i]$

$\forall i: 0 \leq i \leq n: y_i[j] = y_{i-1}[j] + x[j-i] * a[i]$

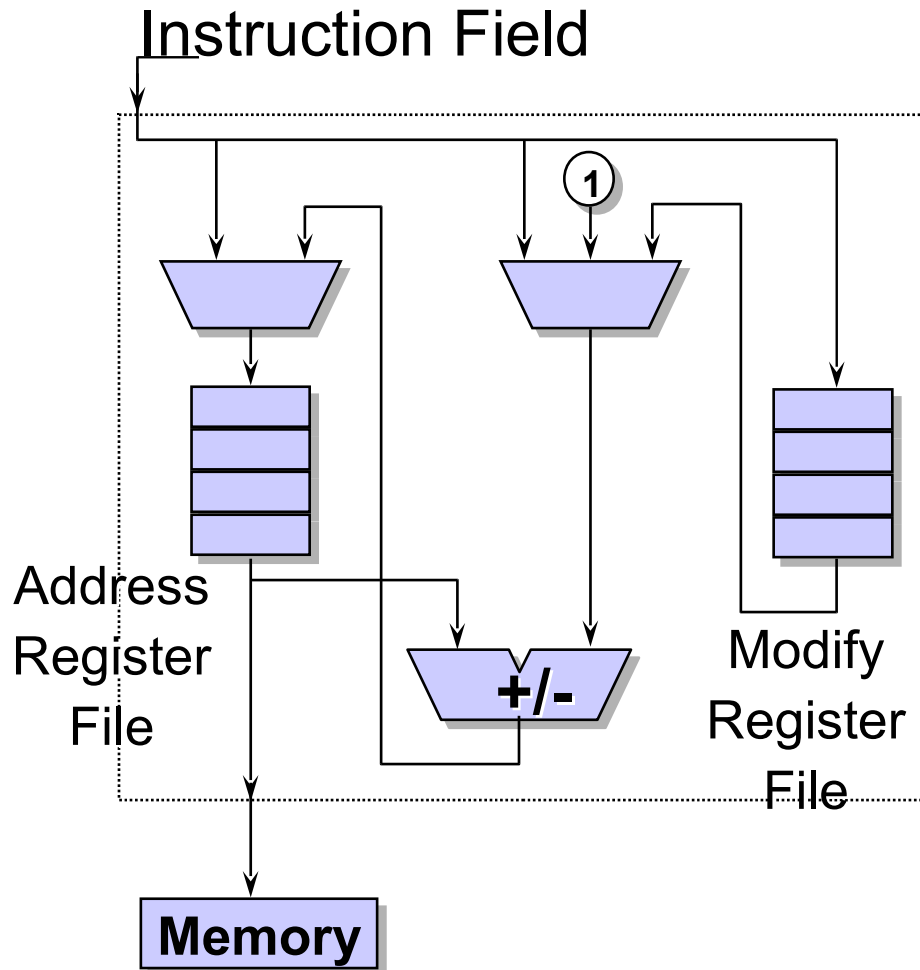
Architecture: Example: Data path ADSP210x



- Parallelism
- Dedicated registers
- No matching compiler \Rightarrow inefficient code

Exploitation of parallel address computations

Generic address generation unit (AGU) model



Parameters:

$k = \#$ address registers

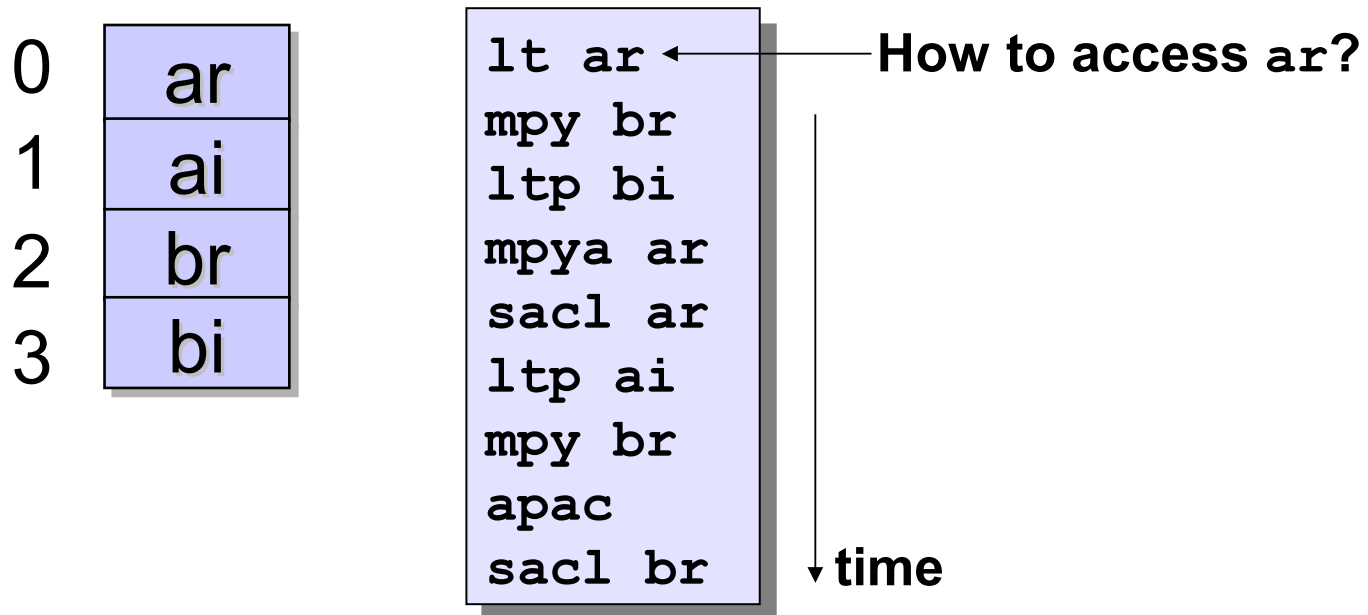
$m = \#$ modify registers

Cost metric for AGU operations:

| Operation | cost |
|------------------------------|------|
| immediate AR load | 1 |
| immediate AR modify | 1 |
| auto-increment/ decrement | 0 |
| AR += MR | 0 |

Address pointer assignment (APA)

Given: Memory layout + assembly code (without address code)



Address pointer assignment (APA) is the sub-problem of finding an allocation of address registers for a given memory layout and a given schedule.

General approach: Minimum Cost Circulation Problem

Let $G = (V, E, u, c)$, with (V, E) : directed graph

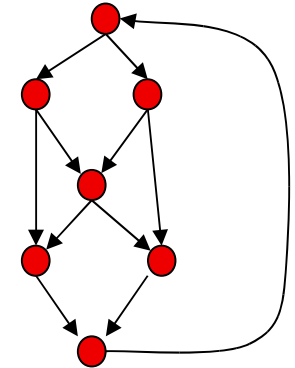
- $u: E \rightarrow \mathbb{R}_{\geq 0}$ is a capacity function,
- $c: E \rightarrow \mathbb{R}$ is a cost function; $n = |V|$, $m = |E|$.

Definition:

1. $g: E \rightarrow \mathbb{R}_{\geq 0}$ is called a **circulation** if it satisfies :

$$\forall v \in V: \sum_{w \in V: (v,w) \in E} g(v,w) = \sum_{w \in V: (w,v) \in E} g(w,v) \quad (\text{flow conservation})$$

2. g is **feasible** if $\forall (v,w) \in E: g(v,w) \leq u(v,w)$ (capacity constraints)
3. The cost of a circulation g is $c(g) = \sum_{(v,w) \in E} c(v,w) g(v,w)$.
4. There may be a lower bound on the flow through an edge.
5. The **minimum cost circulation problem** is to find a feasible circulation of minimum cost.



[K.D. Wayne: A Polynomial Combinatorial Algorithm for Generalized Minimum Cost Flow, <http://www.cs.princeton.edu/~wayne/papers/ratio.pdf>]

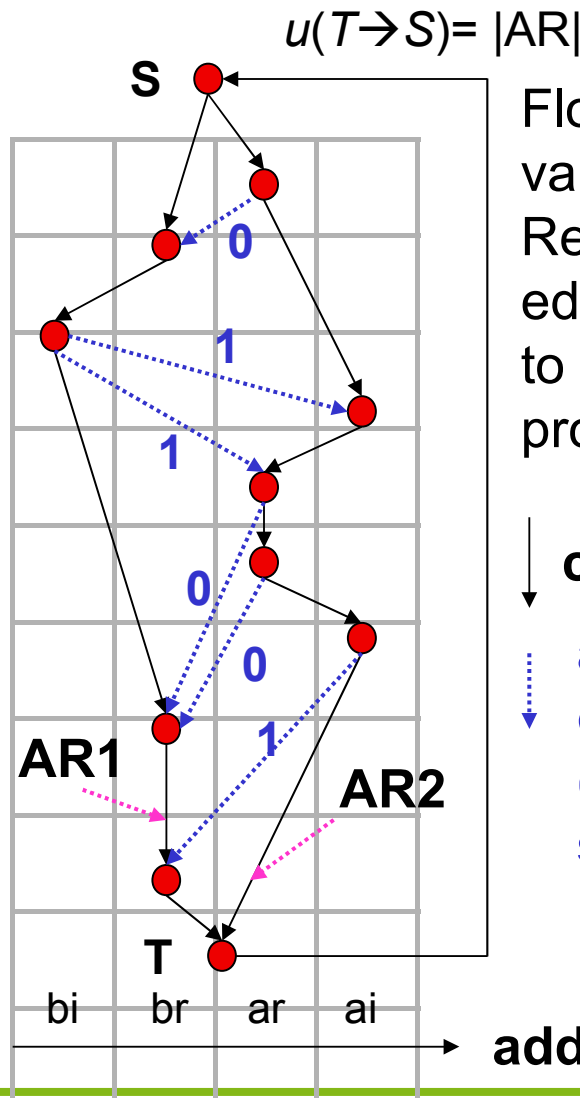
Mapping APA to the Minimum Cost Circulation Problem

Assembly sequence*

```
lt ar
mpy br
ltp bi
mpy ai
mpya ar
sac1 ar
ltp ai
mpy br
apac
sac1 br
```

time

Variables



Flow into and out of variable nodes must be 1. Replace variable nodes by edge with lower bound=1 to obtain pure circulation problem

circulation selected

additional edges of original graph (only samples shown)

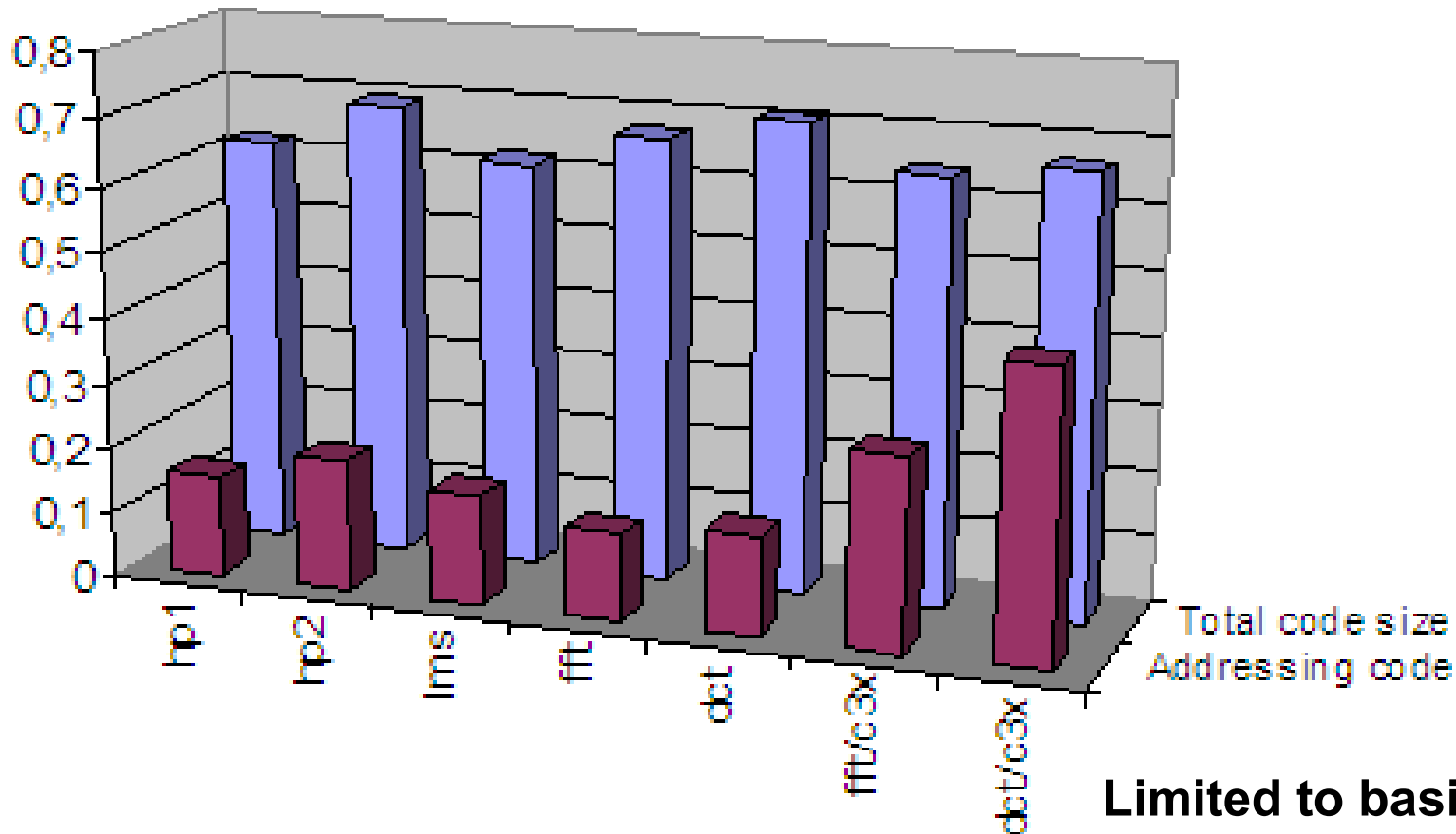
* C2x processor from ti

[C. Gebotys: DSP Address Optimization Using A Minimum Cost Circulation Technique, ICCAD, 1997]

Results according to Gebotys

Optimized code size

Original code size



Limited to basic blocks

Beyond basic blocks: - handling array references in loops -

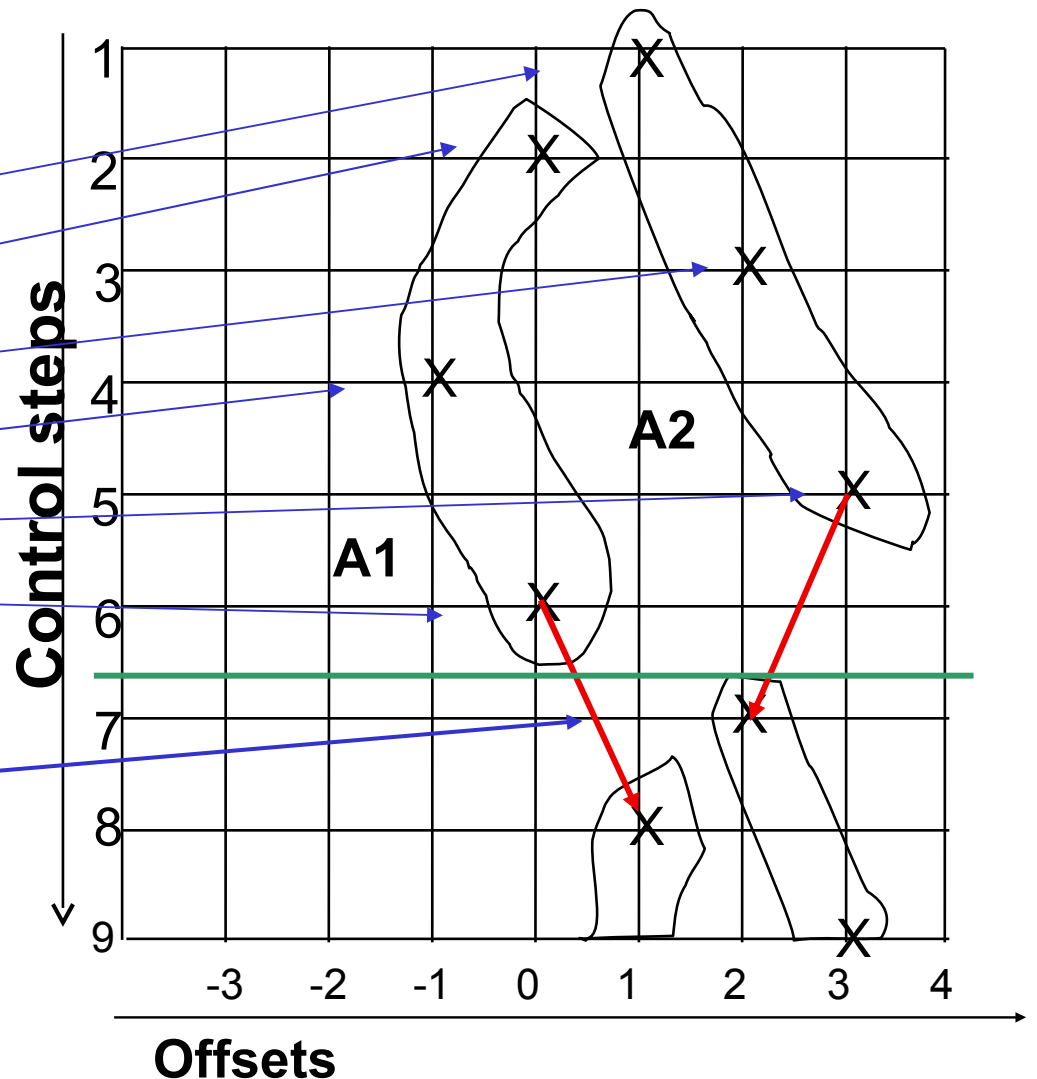
Example:

```

for (i=2; i<=N; i++)
{ .. B[i+1] /*A2++ */
  .. B[i] /*A1-- */
  .. B[i+2] /*A2++ */
  .. B[i-1] /*A1++ */
  .. B[i+3] /*A2-- */
  .. B[i] /*A1++ */
}
    
```

Cost for crossing loop boundaries considered.

Reference: A. Basu, R. Leupers, P. Marwedel: Array Index Allocation under Register Constraints, Int. Conf. on VLSI Design, Goa/India, 1999



Offset assignment problem (OA)

- Effect of optimised memory layout -

Let's assume that we can modify the memory layout

☞ offset assignment problem.

(k,m,r) -OA is the problem of generating a memory layout which minimizes the cost of addressing variables, with

☞ k : number of address registers

☞ m : number of modify registers

☞ r : the offset range

The case $(1,0,1)$ is called simple offset assignment (SOA), the case $(k,0,1)$ is called general offset assignment (GOA).

☞ SOA example

- Effect of optimised memory layout -

Variables in a basic block: Access sequence:

$V = \{a, b, c, d\}$

$S = (b, d, a, c, d, c)$

| | | |
|---|---|--------------|
| 0 | a | Load AR,1 ;b |
| 1 | b | AR += 2 ;d |
| 2 | c | AR -= 3 ;a |
| 3 | d | AR += 2 ;c |
| | | AR ++ ;d |
| | | AR -- ;c |

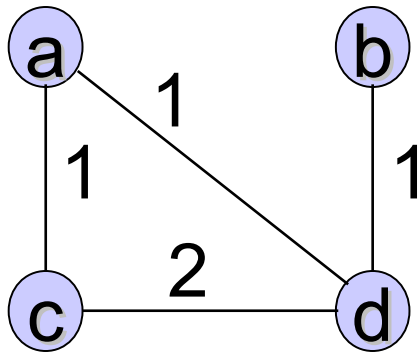
cost: 4

| | | |
|---|---|--------------|
| 0 | b | Load AR,0 ;b |
| 1 | d | AR ++ ;d |
| 2 | c | AR +=2 ;a |
| 3 | a | AR -- ;c |
| | | AR -- ;d |
| | | AR ++ ;c |

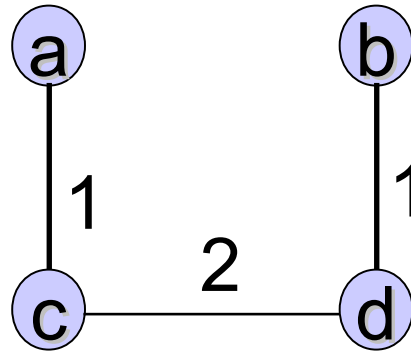
cost: 2

SOA example: Access sequence, access graph and Hamiltonian paths

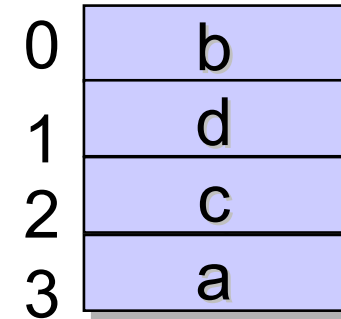
access sequence: b d a c d c



access graph



maximum weighted path=
max. weighted Hamiltonian
path covering (MWHC)



memory layout

SOA used as a building block for more complex situations

➔ significant interest in good SOA algorithms

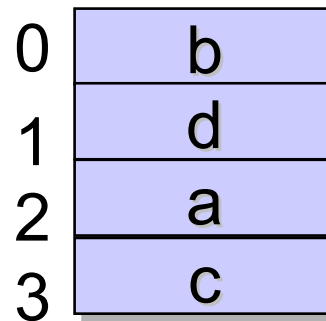
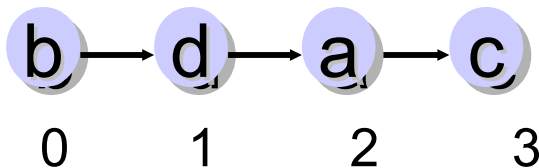
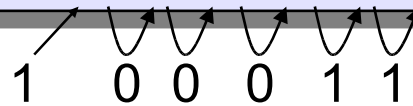
[Bartley, 1992; Liao, 1995]

Naïve SOA

Nodes are added in the order in which they are used in the program.

Example:

Access sequence: $S = (b, d, a, c, d, c)$



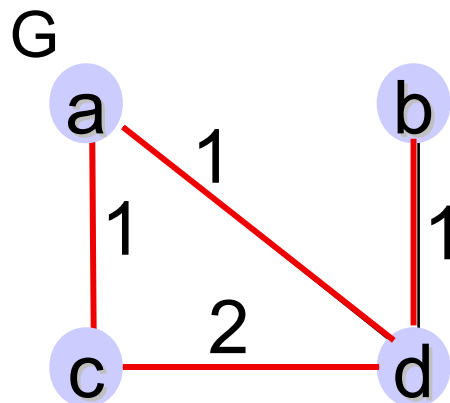
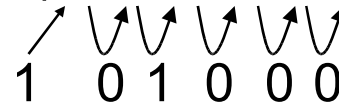
memory layout

Liao's algorithm

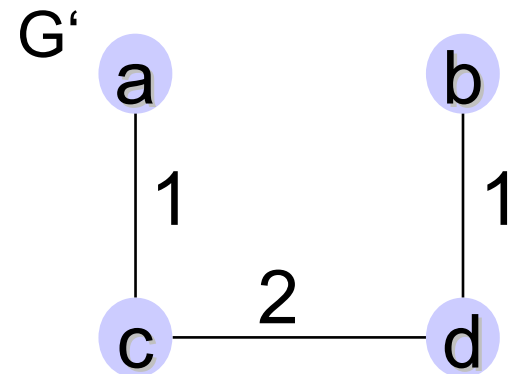
Similar to Kruskal's spanning tree algorithms:

1. Sort edges of access graph $G=(V,E)$ according to their weight
2. Construct a new graph $G'=(V',E')$, starting with $E'=0$
3. Select an edge e of G of highest weight; If this edge does not cause a cycle in G' and does not cause any node in G' to have a degree > 2 then add this node to E' otherwise discard e .
4. Goto 3 as long as not all edges from G have been selected and as long as G' has less than the maximum number of edges $(|V|-1)$.

Example: Access sequence: $S=(b, d, a, c, d, c)$



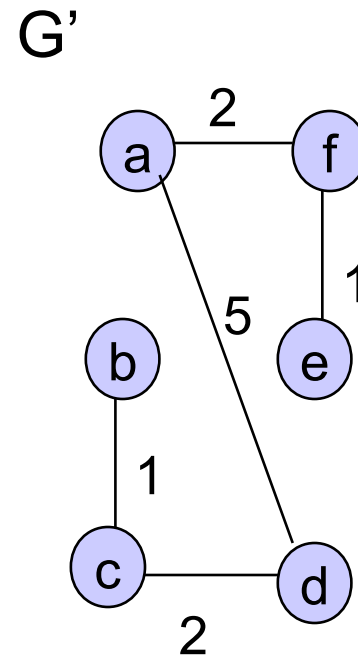
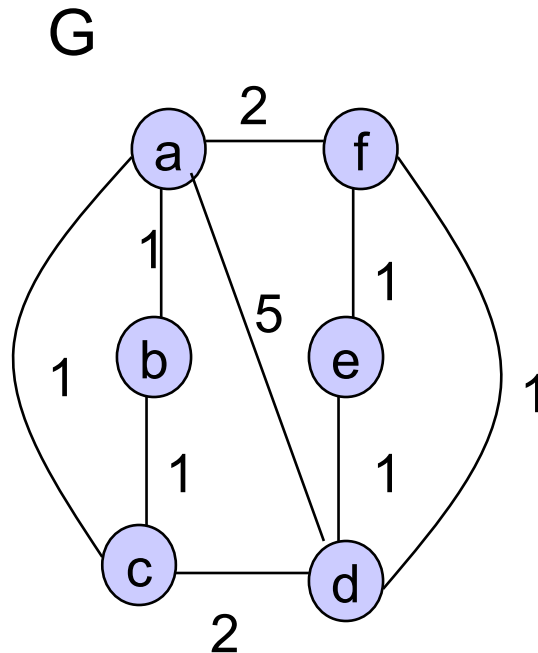
- 2 (c,d)
- 1 (a,c)
- 1 (a,d)
- 1 (b,d)



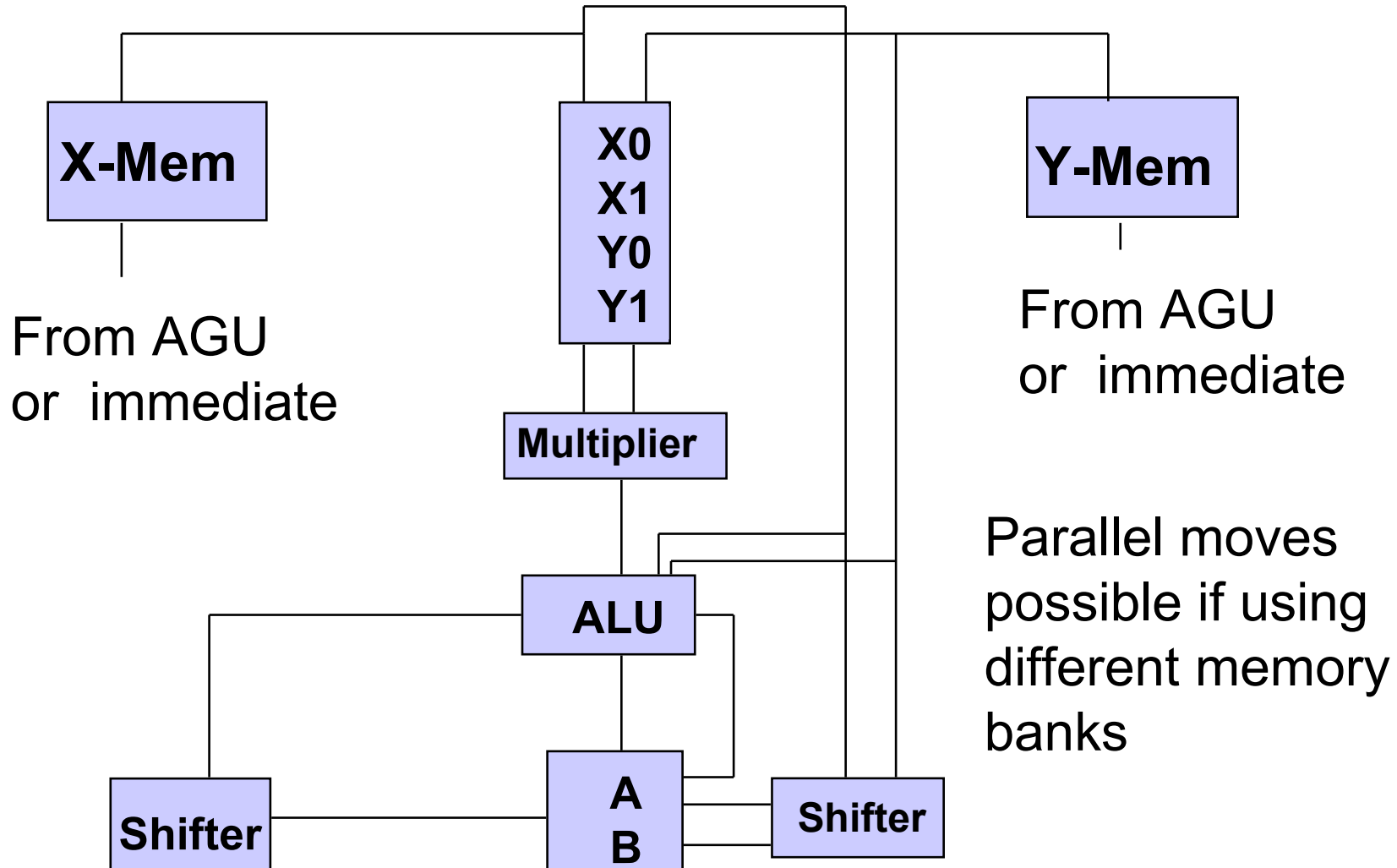
Implicit edges of weight 0 for all unconnected nodes.

Liao's algorithm on a more complex graph

a b c d e f a d a d a c d f a d



Multiple memory banks - Sample hardware -



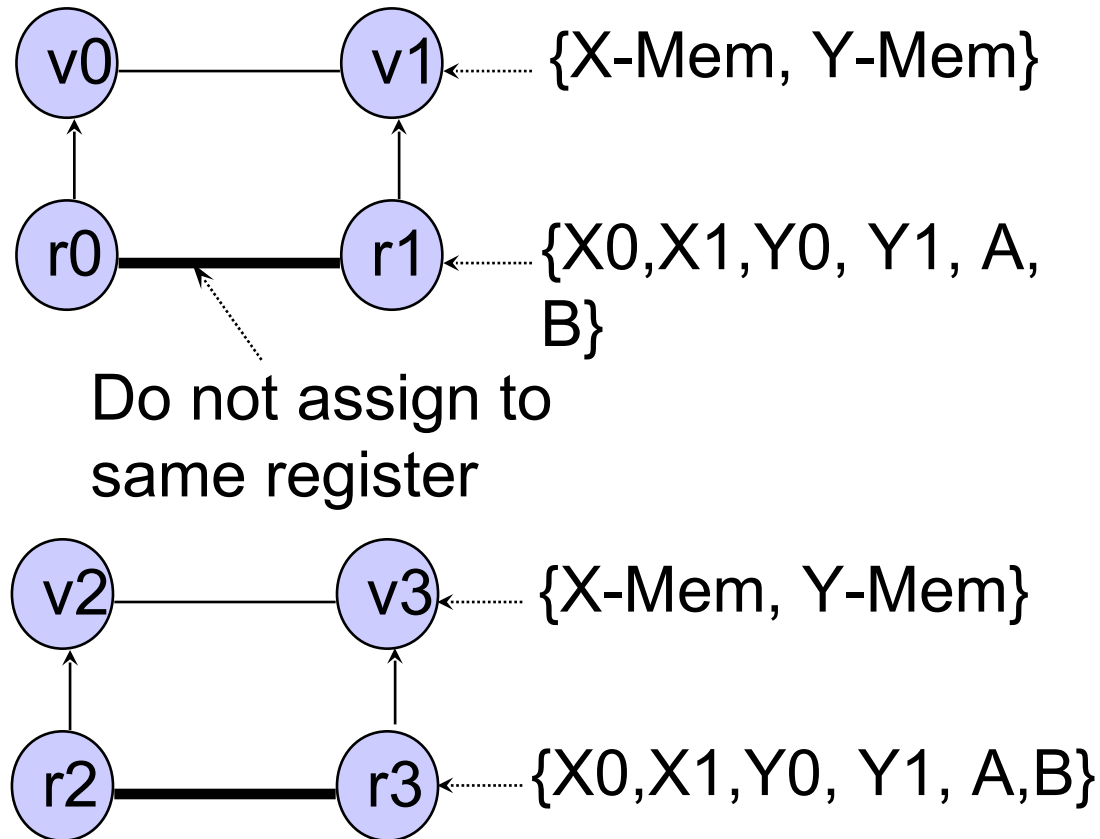
Multiple memory banks

- Constraint graph generation -

Precompacted
code
(symbolic variables
and registers)

Move v0,r0 v1,r1
Move v2,r2 v3,r3

Constraint graph



Links maintained, more constraints ...

Further optimizations

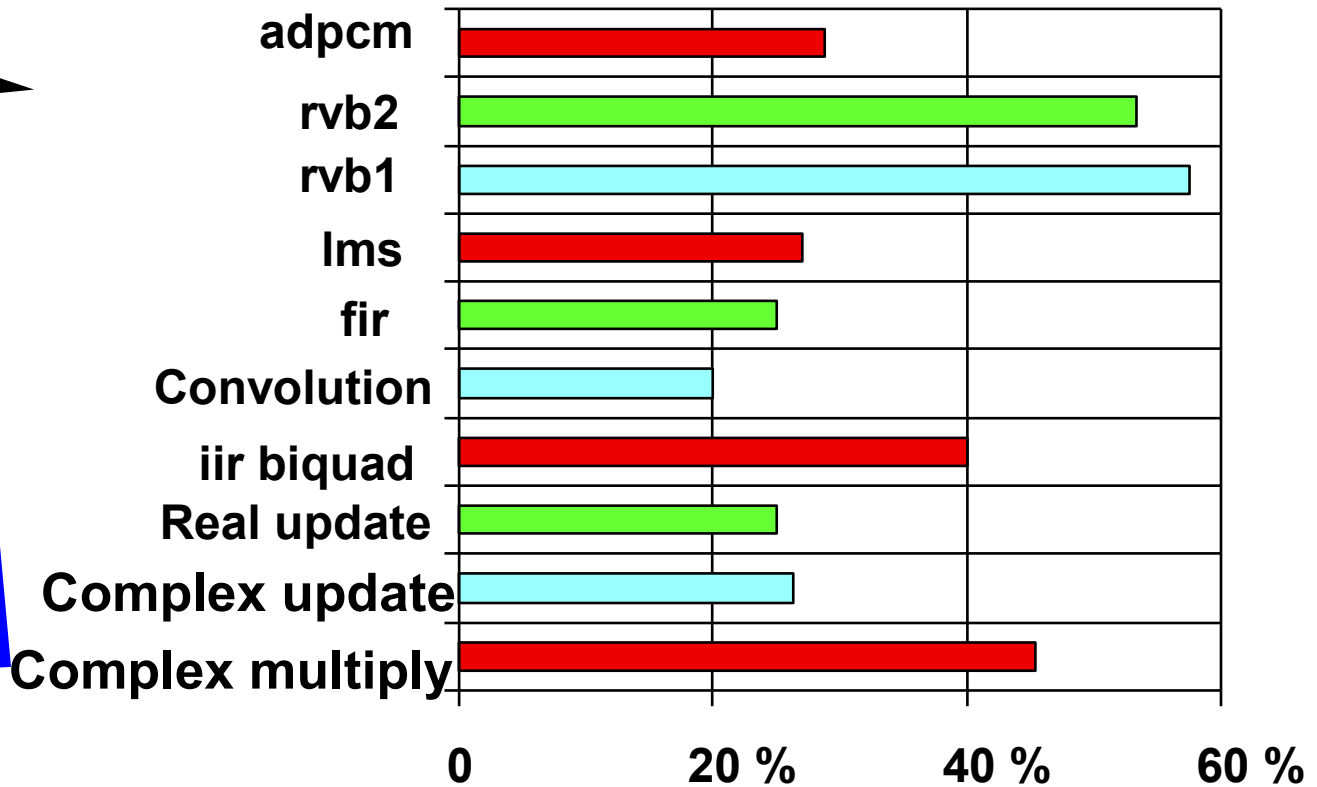
Peter Marwedel
TU Dortmund
Informatik 12
Germany

2010/01/13



Multiple memory banks

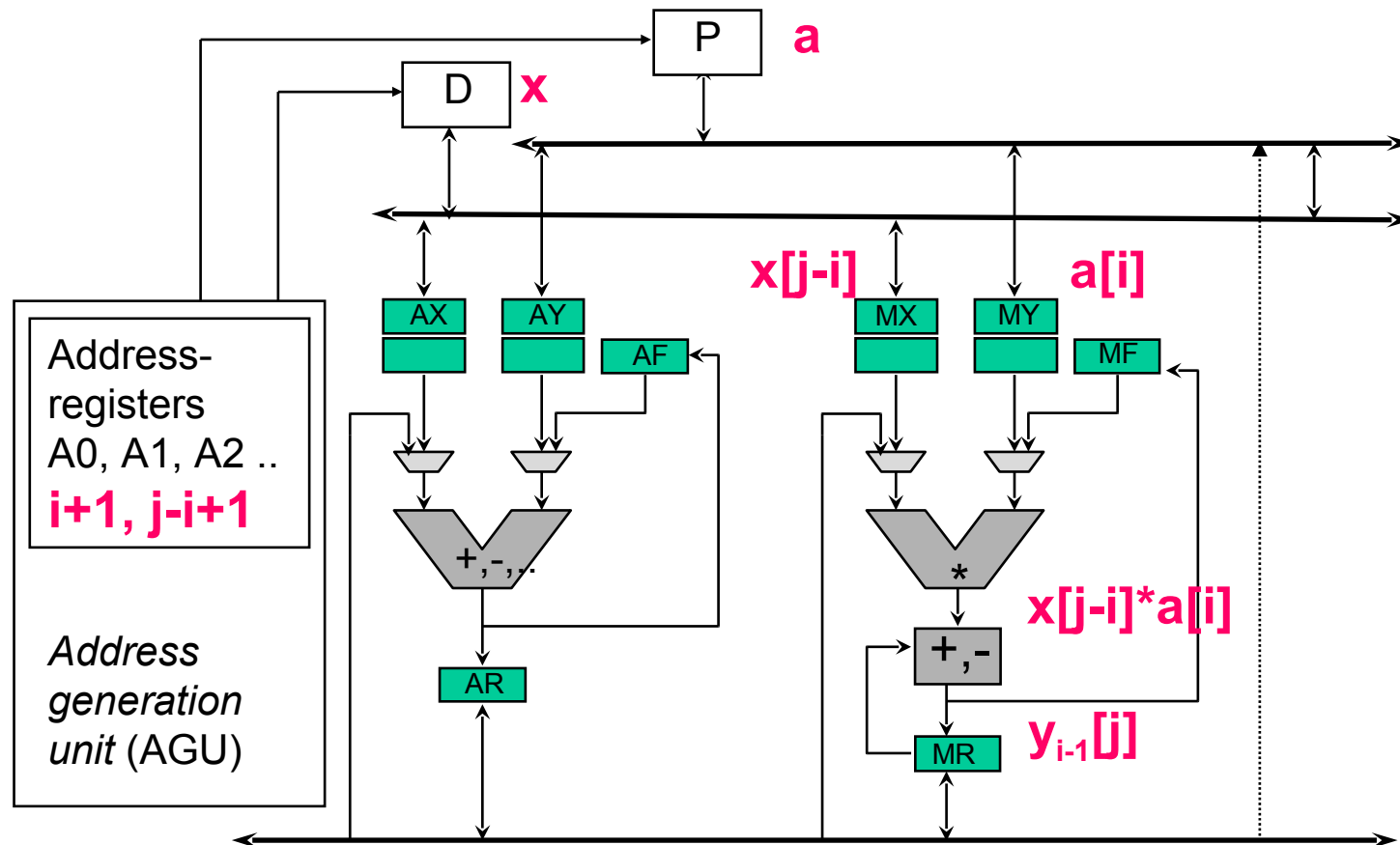
Code size reduction through simulated annealing



[Sudarsanam, Malik, 1995]

Exploitation of instruction level parallelism (ILP)

Several transfers in the same cycle:



Exploitation of instruction level parallelism (ILP)

1: MR := MR+(MX*MY);

2: MX:=D[A1];

3: MY:=P[A2];

4: A1- -;

5: A2++;

6: D[0]:= MR;

.....



1': MR := MR+(MX*MY), MX:=D[A1],
MY:=P[A2], A1- -, A2++;

2': D[0]:= MR;

Modelling of possible parallelism using n-ary compatibility relation, e.g. $\sim(1,2,3,4,5)$

Generation of integer programming (IP)- model (max. 50 statements/model)

Using standard-IP-solver to solve equations

Exploitation of instruction level parallelism (ILP)

$$u(n) = u(n - 1) + K0 \times e(n) + K1 \times e(n - 1);$$

$$e(n - 1) = e(n)$$



$$\text{ACCU} := u(n - 1)$$

$$\text{TR} := e(n - 1)$$

$$\text{PR} := \text{TR} \times K1$$

$$\text{TR} := e(n)$$

$$e(n - 1) := e(n)$$

$$\text{ACCU} := \text{ACCU} + \text{PR}$$

$$\text{PR} := \text{TR} \times K0$$

$$\text{ACCU} := \text{ACCU} + \text{PR}$$

$$u(n) := \text{ACCU}$$

$$\text{ACCU} := u(n - 1)$$

$$\text{TR} := e(n - 1)$$

$$\text{PR} := \text{TR} \times K1$$

$$e(n - 1) := e(n) \parallel \text{TR} := e(n) \parallel$$

$$\text{ACCU} := \text{ACCU} + \text{PR}$$

$$\text{PR} := \text{TR} \times K0$$

$$\text{ACCU} := \text{ACCU} + \text{PR}$$

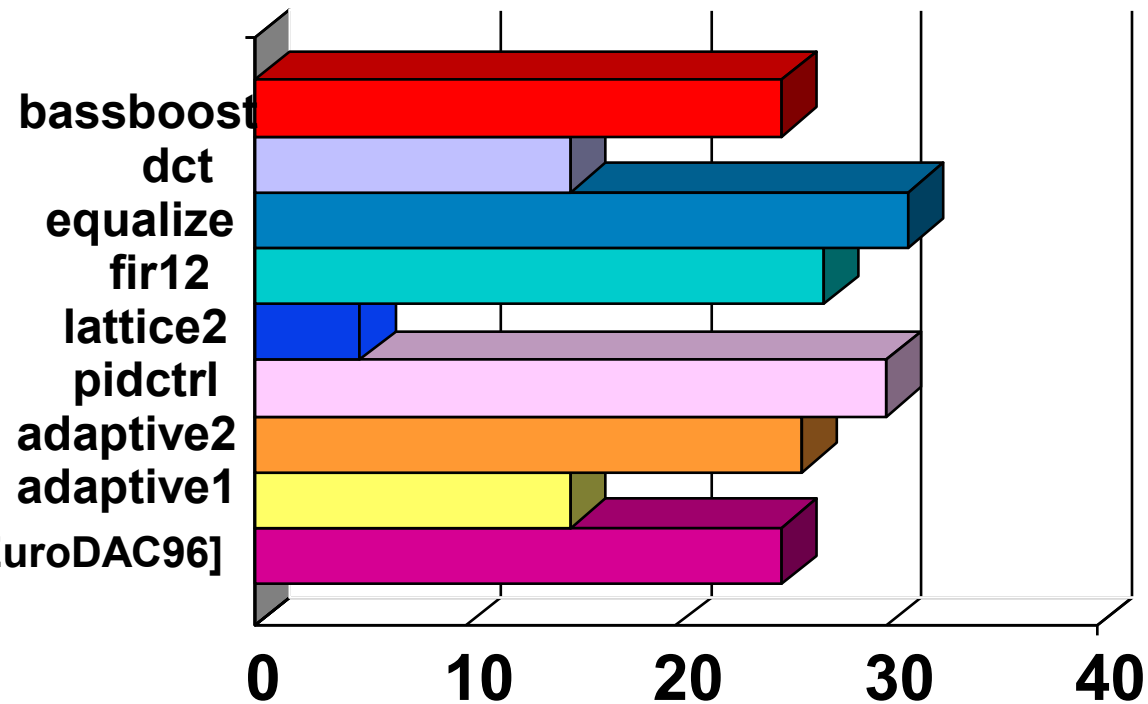
$$u(n) := \text{ACCU}$$

- From 9 to 7
cycles
through
compaction

Exploitation of instruction level parallelism (ILP)

Results obtained through integer programming:

Code size reduction [%]



[Leupers, EuroDAC96]

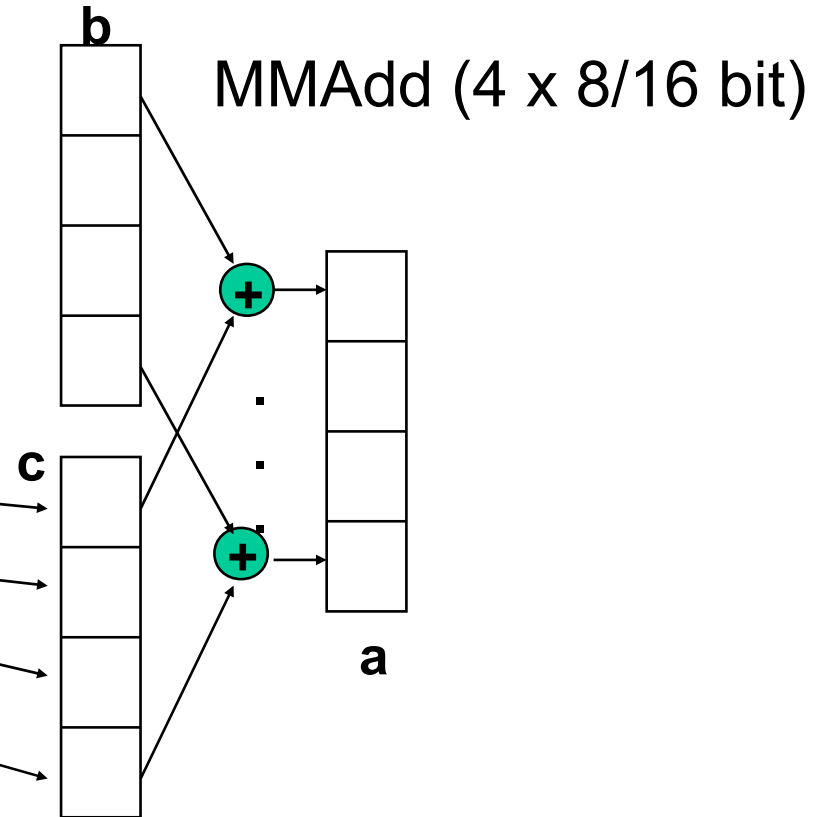
Compaction times: 2 .. 35 sec

Exploitation of Multimedia Instructions

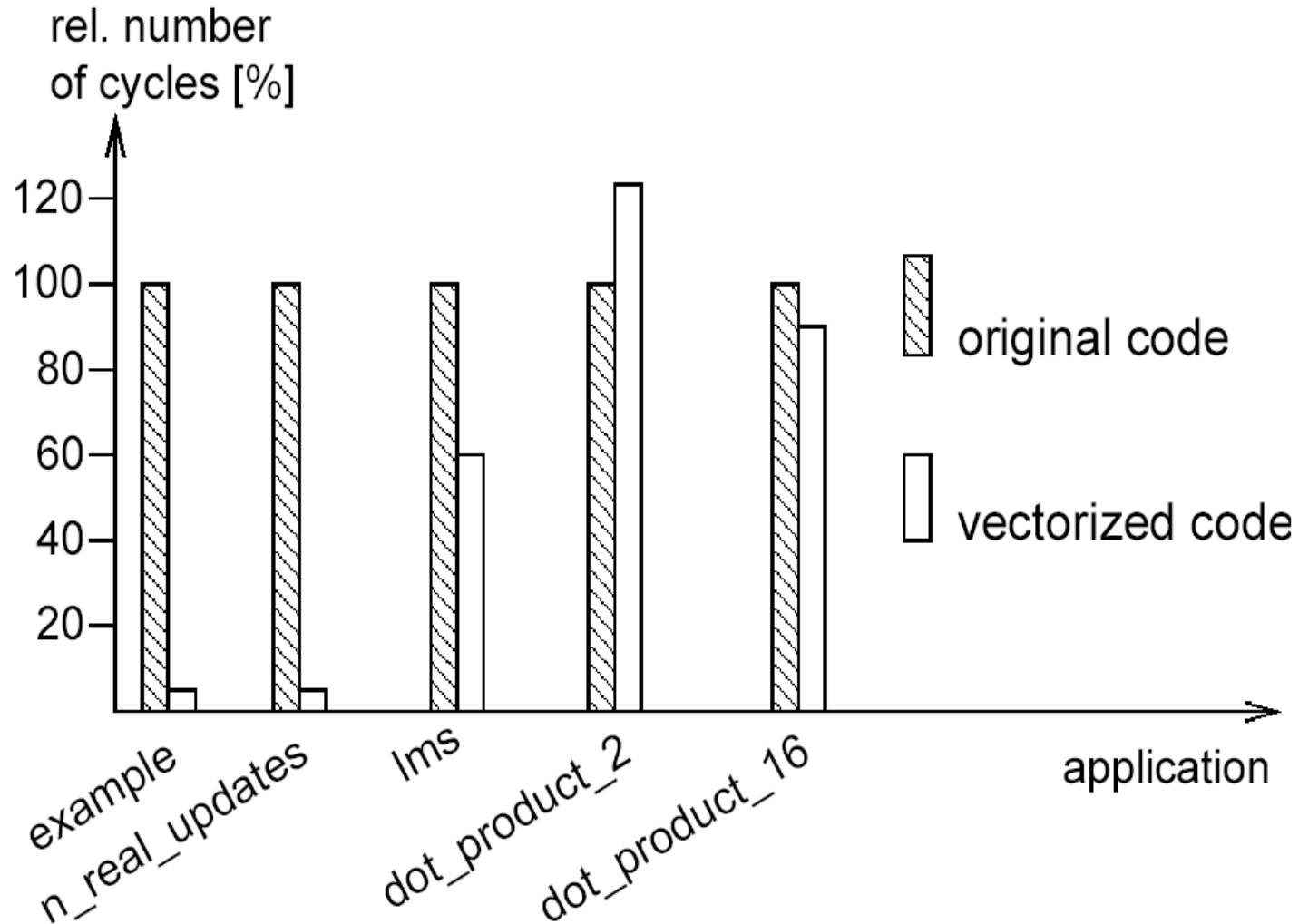
```
FOR i:=0 TO n DO  
  a[i] = b[i] + c[i]
```



```
FOR i:=0 STEP 4 TO n DO  
  a[i ]=b[i ]+c[i ] ;  
  a[i+1]=b[i+1]+c[i+1] ;  
  a[i+2]=b[i+2]+c[i+2] ;  
  a[i+3]=b[i+3]+c[i+3] ;
```



Improvements for M3 DSP due to vectorization

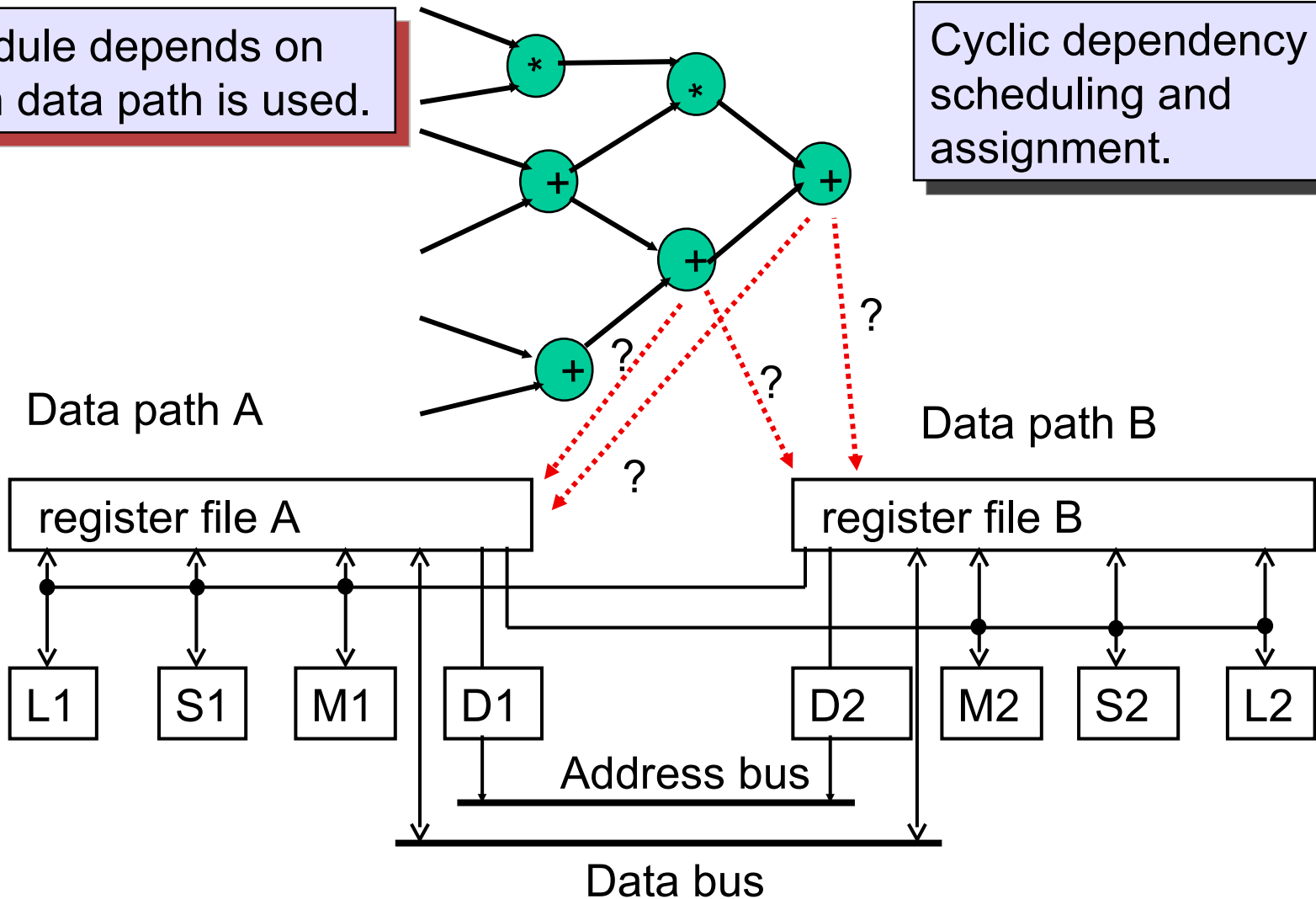


Scheduling for partitioned data paths

Schedule depends on which data path is used.

Cyclic dependency of scheduling and assignment.

'C6x:

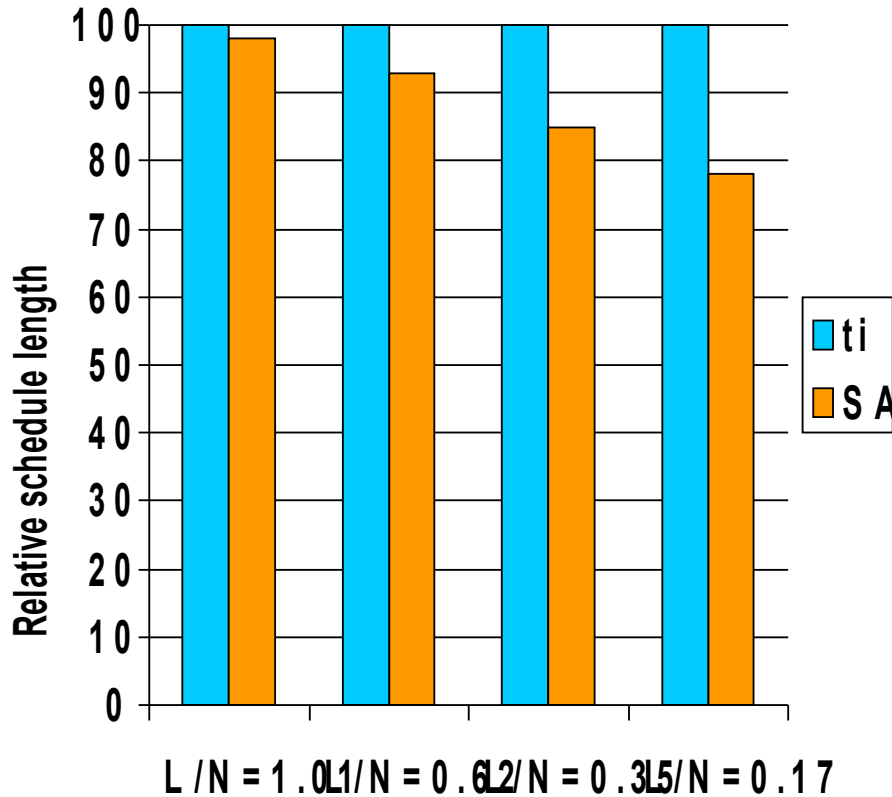


Integrated scheduling and assignment using Simulated Annealing (SA)

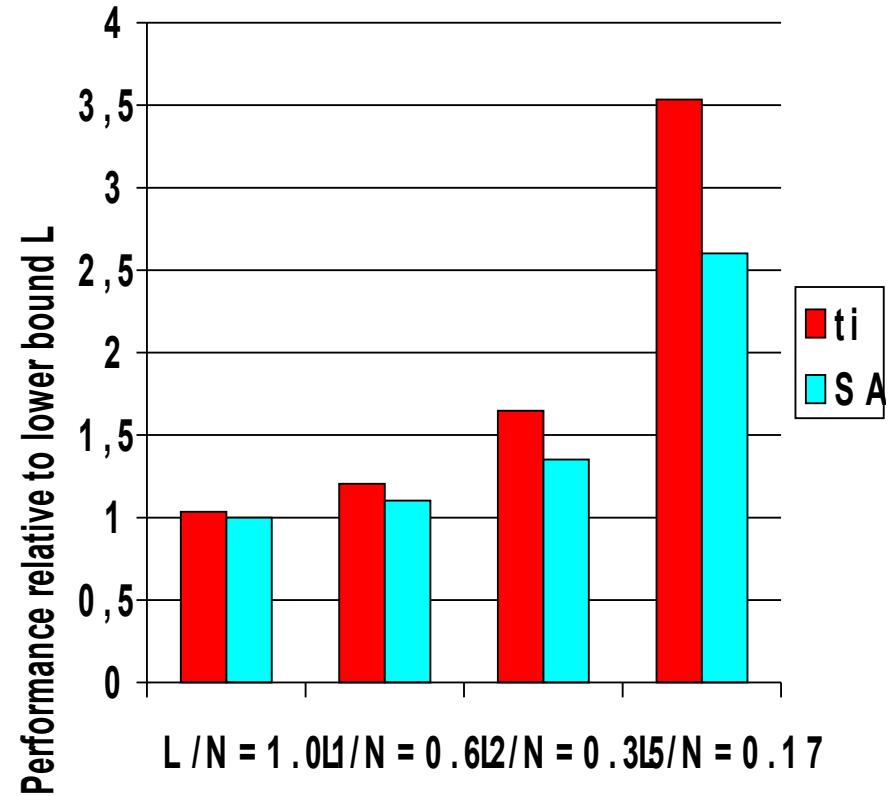
```
algorithm Partition
input DFG G with nodes;
output: DP: array [1..N] of 0,1 ;
var int i, r, cost, mincost;
float T;
begin
  T=10;
  DP:=Randompartitioning;
  mincost :=
    LISTSCHEDELING(G,D,P);
  WHILE_LOOP;
  return DP;
end.
```

```
WHILE_LOOP:
while T>0.01 do
  for i=1 to 50 do
    r:= RANDOM(1,n);
    DP[r] := 1-DP[r];
    cost:=LISTSCHEDELING(G,D,P);
    delta:=cost-mincost;
    if delta <0 or
      RANDOM(0,1)<exp(-delta/T)
    then mincost:=cost
    else DP[r]:=1-DP[r]
    end if;
  end for;
  T:= 0.9 * T;
end while;
```

Results: relative schedule length as a function of the “width” of the DFG



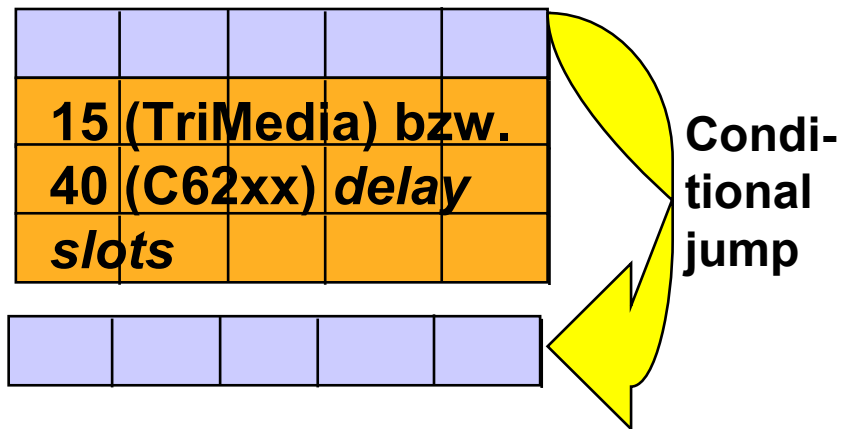
SA approach outperforms the ti approach for “wide” DFGs (containing a lot of parallelism)



For wide DFGs, SA algorithm is able of “staying closer” critical path length.

VLIW (very long instruction word) DSPs

Large *branch delay penalty*:



Avoiding this penalty:
predicated execution:

[c] *instruction*

c=*true*: instruction executed

c=*false*: effectively NOOP

Realisation of *if-statements*

with conditional jumps or with *predicated execution*:

```

if (c)
{ a = x + y;
  b = x + z;
}
else
{ a = x - y;
  b = x - z;
}
    
```

Cond. instructions:

```

[c] ADD x,y,a
|| [c] ADD x,z,b
|| [!c] SUB x,y,a
|| [!c] SUB x,z,b
    
```

1 cycle

Cost of implementation methods for IF-Statements

Sourcecode: if (c1) {t1; if (c2) t2}

No precondition (no enclosing IF or enclosing IFs implemented with cond. jumps)

1. Conditional jump:

BNE c1, L;

t1;

L: ...

2. Conditional

Instruction:

[c1] t1

Precondition (enclosing IF not implemented with conditional jumps)

3. Conditional jump :

[c1] c:=c2

[~c1] c:=0

BNE c, L;

t2;

L: ...

4. Conditional

Instruction :

[c1] c:=c2

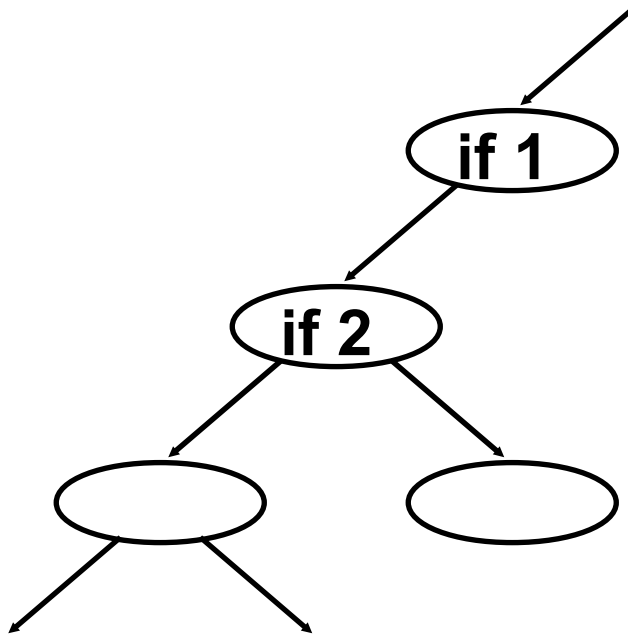
[~c1] c:=0

[c] t2

Additional computations to compute effective condition c

Optimization for nested IF-statements

Goal: compute fastest implementation for all IF-statements



- Selection of fastest implementation for if-1 requires knowledge of how fast if-2 can be implemented.
- Execution time of if-2 depends on setup code, and, hence, also on how if 1 is implemented
- cyclic dependency!

Dynamic programming algorithm (phase 1)

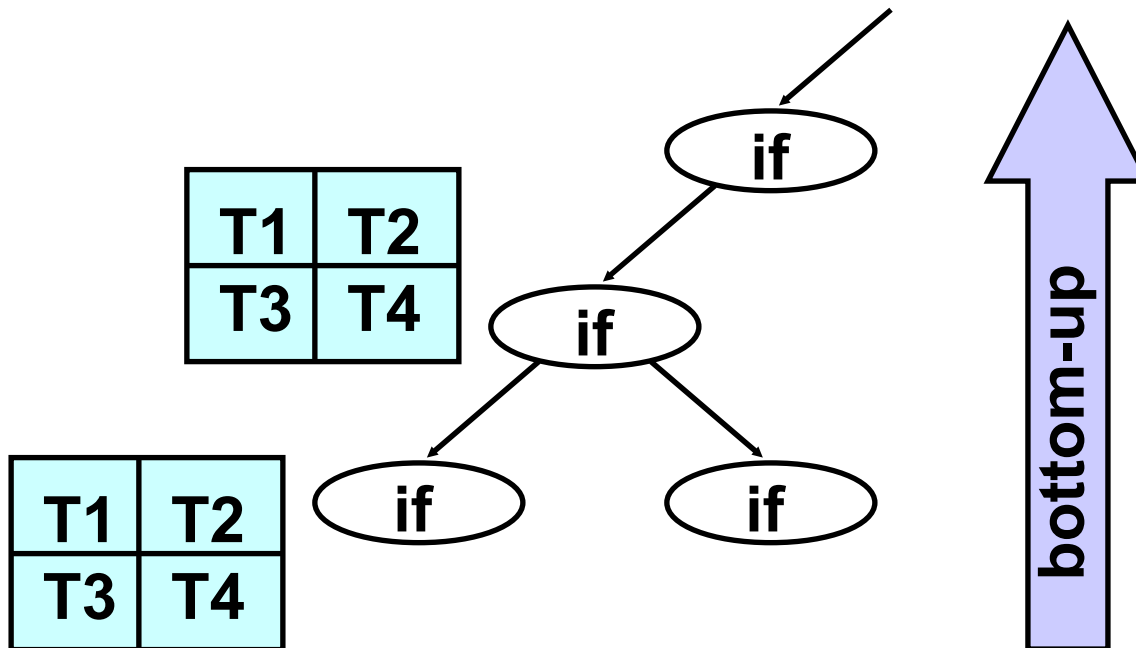
For each if-statement compute 4 cost values:

T1 : cond. jump, no precondition

T2 : cond. instructions, no precondition

T3 : cond. jump, with precondition

T4: cond. instructions, with precondition



Dynamic programming (phase 2)

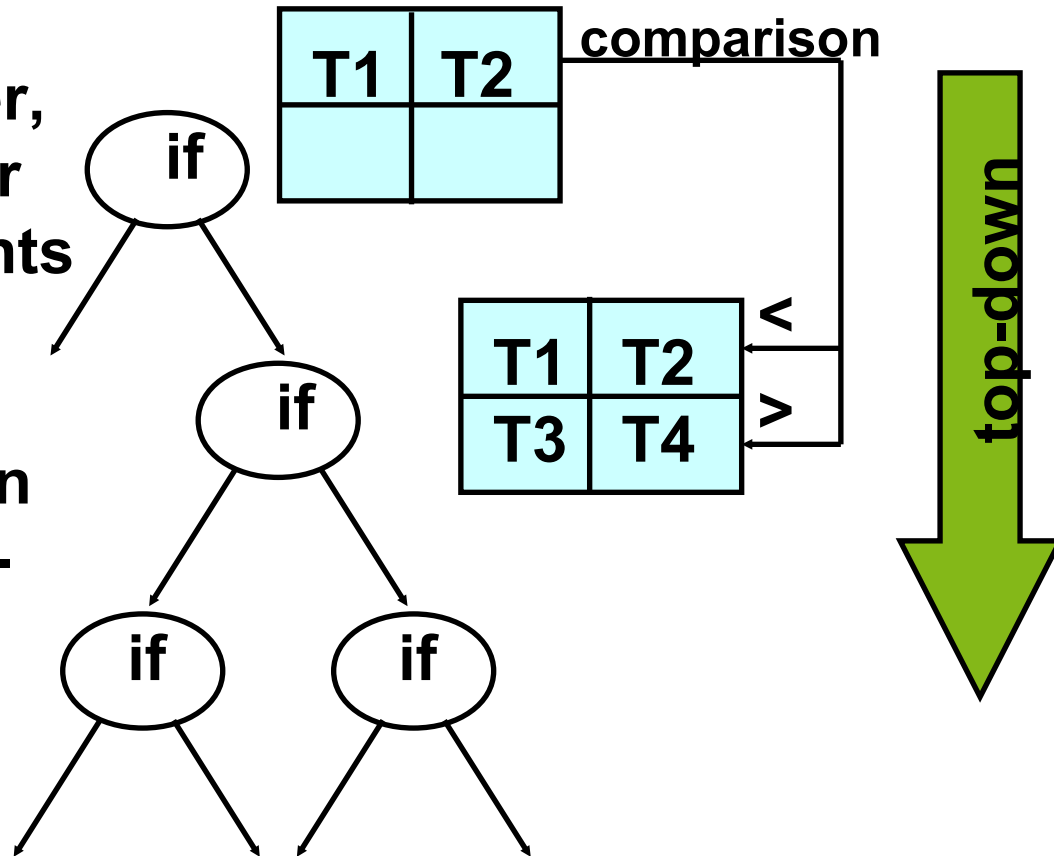
No precondition for top-level IF-statement.
Hence, comparison „ $T1 < T2$ “ suffices.

$T1 < T2$:

cond. branch faster,
no precondition for
nested IF-statements

$T1 > T2$:

cond. instructions
faster, precondition
for nested IF-state-
ments



Results: TI C62xx

Runtimes (max) for 10 control-dominated examples

| Example | Conditional jumps | Conditional instructions | Dynamic program. | Min (col. 2-5) | TI C compiler |
|---------|-------------------|--------------------------|------------------|----------------|---------------|
| 1 | 21 | 11 | 11 | 11 | 15 |
| 2 | 12 | 13 | 13 | 12 | 13 |
| 3 | 26 | 21 | 22 | 21 | 27 |
| 4 | 9 | 12 | 12 | 9 | 10 |
| 5 | 26 | 30 | 24 | 24 | 21 |
| 6 | 32 | 23 | 23 | 23 | 30 |
| 7 | 57 | 173 | 49 | 49 | 51 |
| 8 | 39 | 244 | 30 | 30 | 41 |
| 9 | 28 | 27 | 27 | 27 | 29 |
| 10 | 27 | 30 | 30 | 27 | 28 |

Average gain: 12%

Function inlining: advantages and limitations

Advantage: low calling overhead

```
Function sq(c:integer)
    push PC;
    push b;
    BRA sq;
    pull R1;
    mul R1,R1,R1;
    pull R2;
    push R1;
    BRA (R2)+1;
    pull R1;
    ST R1,a;
end;

return:integer;
begin
    return c*c
end;

....
a=sq(b);
....
```

branching

```
....
LD R1,b;
MUL R1,R1,R1;
ST R1,a
```

Inlining

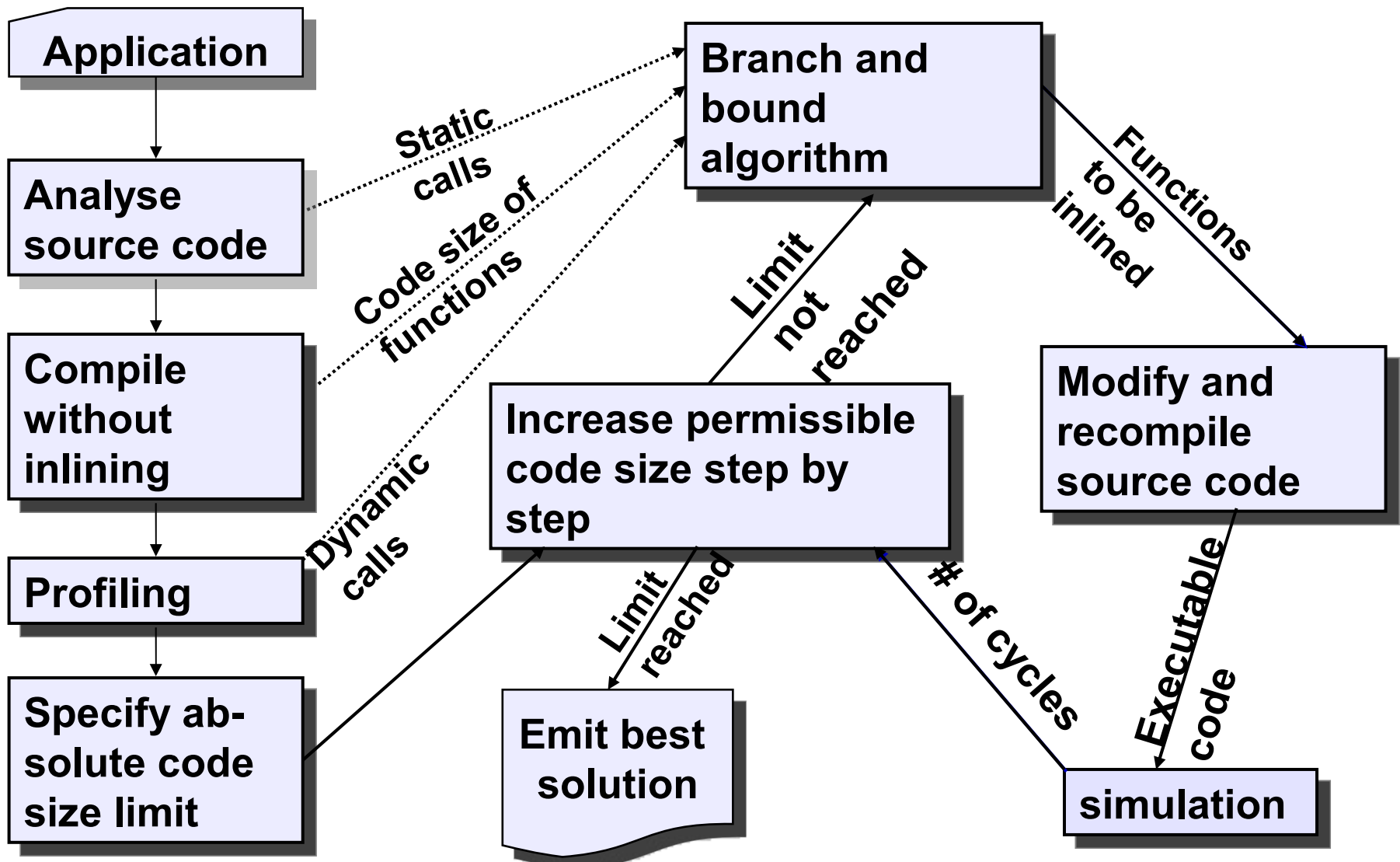
Limitations:

- Not all functions are candidates.
- Code size explosion.
- Requires manual identification using 'inline' qualifier.

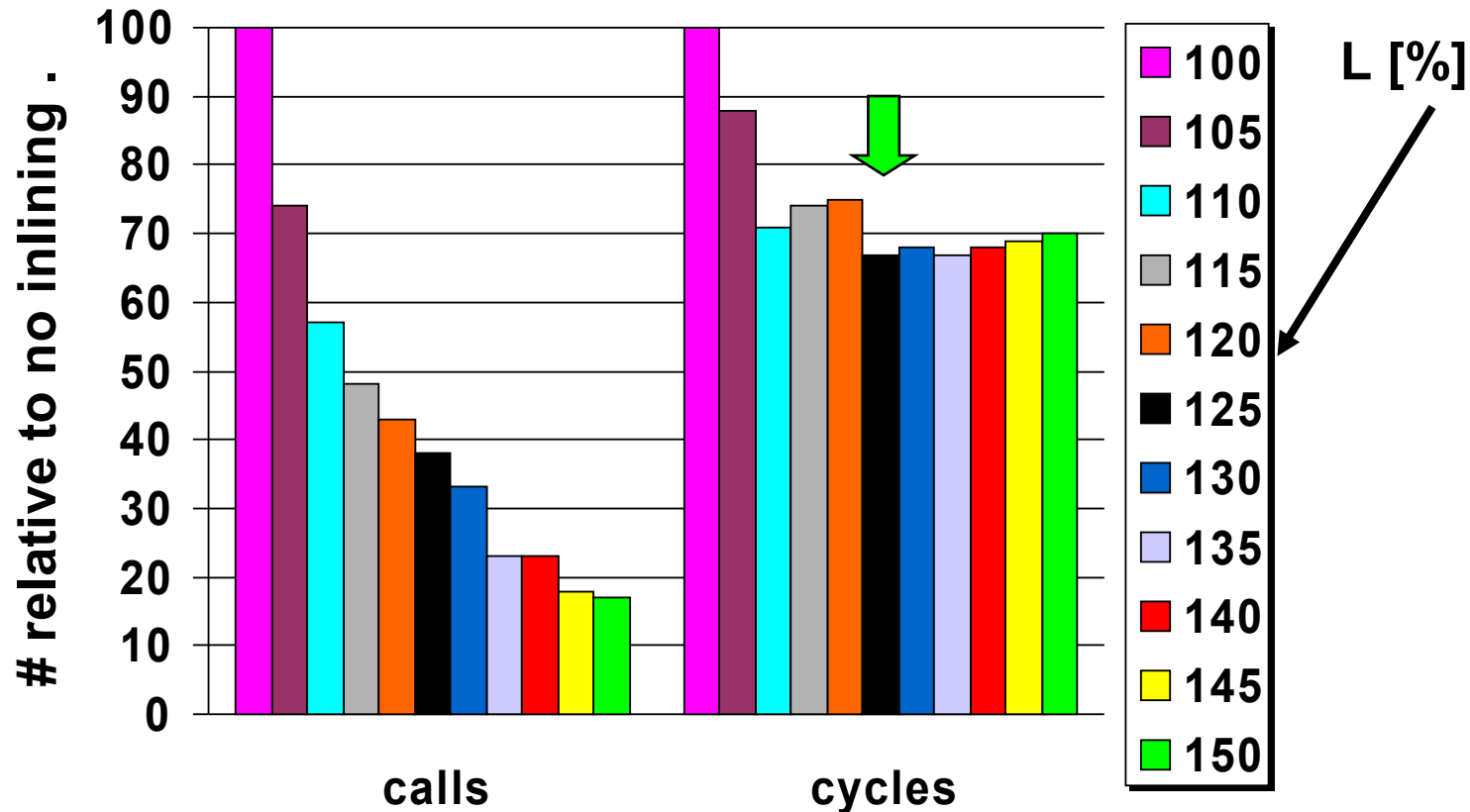
Goal:

- Controlled code size
- Automatic identification of suitable functions.

Design flow



Results for GSM speech and channel encoder: #calls, #cycles (TI 'C62xx)



33% speedup for 25% increase in code size.

of cycles not a monotonically decreasing function of the code size!

Inline vectors computed by B&B algorithm

| size limit (%) | inline vector (functions 1-26) |
|----------------|----------------------------------|
| 100 | 00000000000000000000000000000000 |
| 105 | 001000000001100001111011111111 |
| 110 | 101110010111100001111111111111 |
| 115 | 101100000000001001000111001 |
| 120 | 101101001010001001101111101 |
| 125 | 101100000001010000100111101 |
| 130 | 001100000000010100100111000 |
| 135 | 101100100011110101110111101 |
| 140 | 101110111111111010111111111111 |
| 145 | 1011011010101010100110111101 |
| 150 | 1011011000001010110110111101 |

Major changes for each new size limit. Difficult to generate manually.

References:

- J. Teich, E. Zitzler, S.S. Bhattacharyya. 3D Exploration of Software Schedules for DSP Algorithms, CODES'99
- R. Leupers, P. Marwedel: Function Inlining under Code Size Constraints for Embedded Processors ICCAD, 1999

Summary

- Optimizations for Caches
 - Code Layout transformations
 - Way prediction
- The Offset assignment problem
 - Address pointer assignment problem
 - Simple offset assignment problem
 - General offset assignment problem
- Further optimizations
 - Compaction
 - Multimedia- and VLIW architecture support
 - Predicated execution support
 - Space-aware inlining