technische universität
dortmund

fakultät für informatik
informatik 12

# Optimizations
## - Compilation for Embedded Processors -

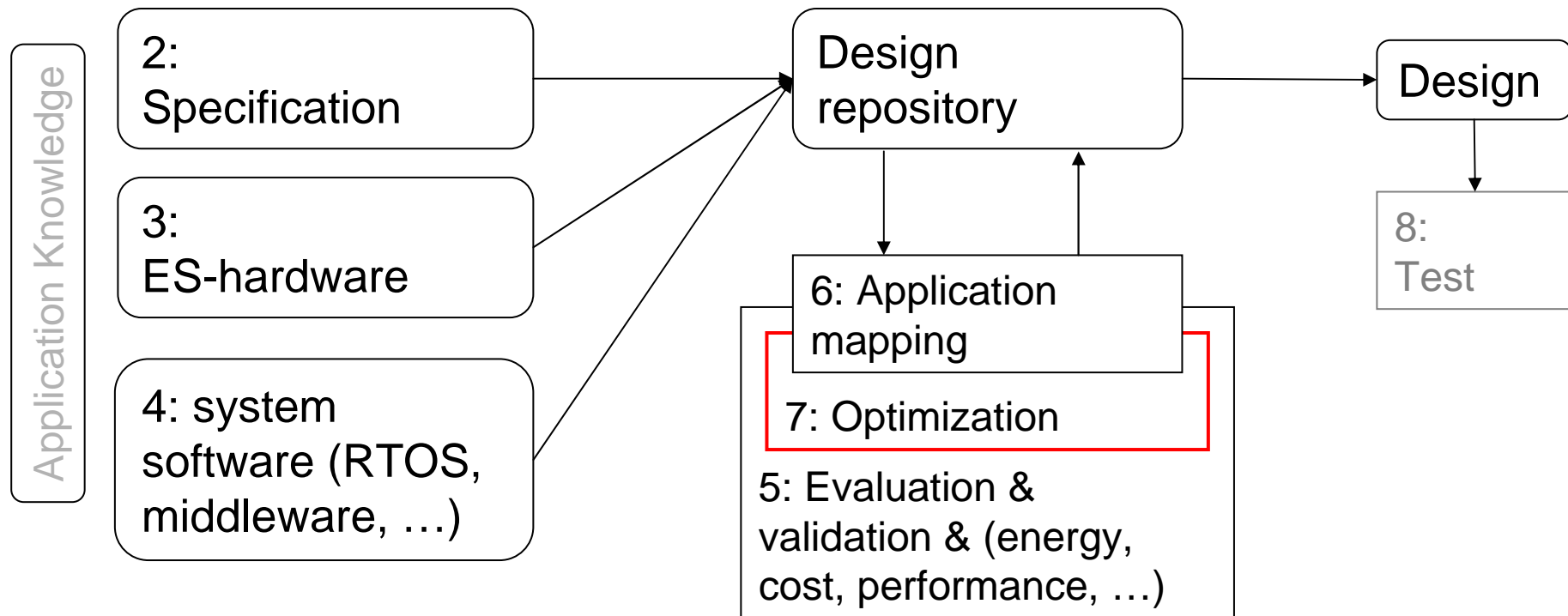Peter Marwedel
TU Dortmund
Informatik 12
Germany

2011年 01 月 25 日

# Structure of this course

Application Knowledge

| 2: Specification |
| 3: ES-hardware |
| 4: system software (RTOS, middleware, …) |

Design repository

Design

8: Test

6: Application mapping

7: Optimization

5: Evaluation & validation & (energy, cost, performance, …)

Numbers denote sequence of chapters
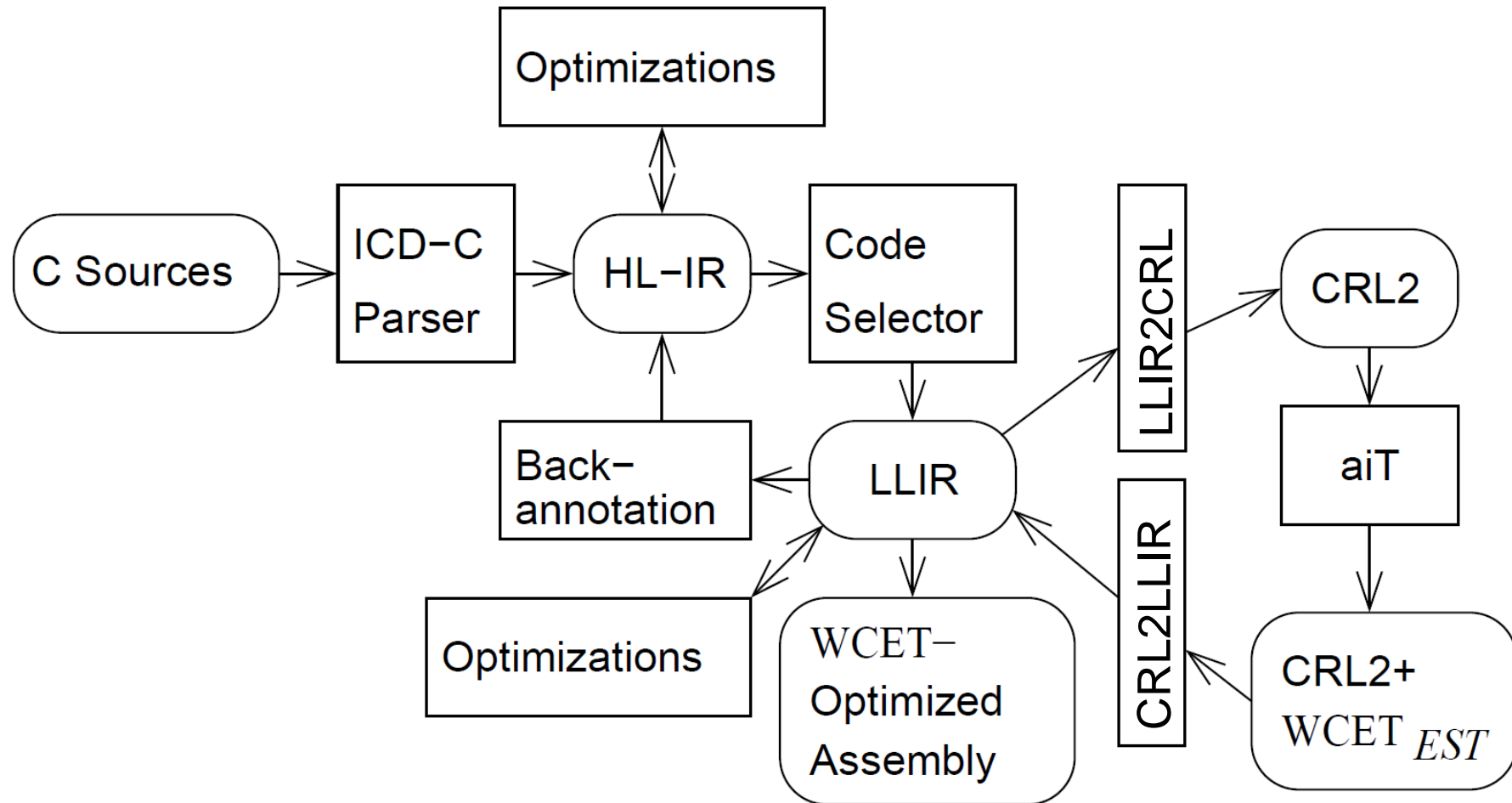
# Reconciling compilers and timing analysis

Compilers mostly unaware of execution times

- Execution times are checked in a "trial-and-error" loop:
  {try: compile – run – check – error: change}*

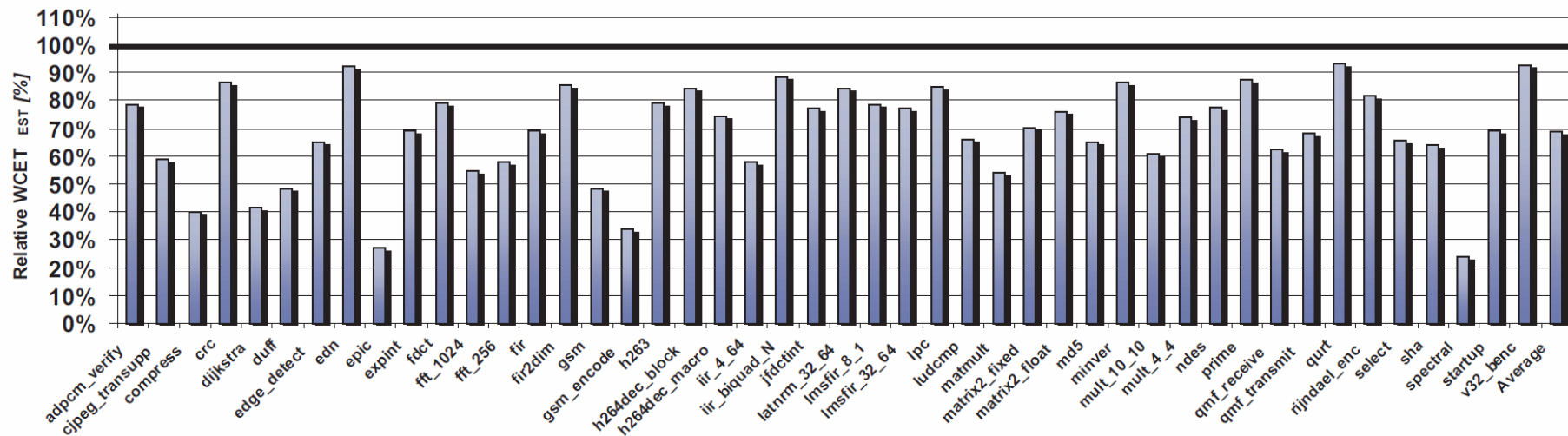- Timing analysis based on measurements is not safe

☞ Integration of safe, static timing analysis into compiler

- Getting rid of loops (if everything works well)

- Safe timing verification

- Potential for optimizing for the WCET

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 3 -

# Structure of WCC (WCET-aware C-compiler)

# Results for WCET-aware register allocation

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 5 -

# The offset assignment problem

Peter Marwedel
TU Dortmund
Informatik 12
Germany

2010/01/13

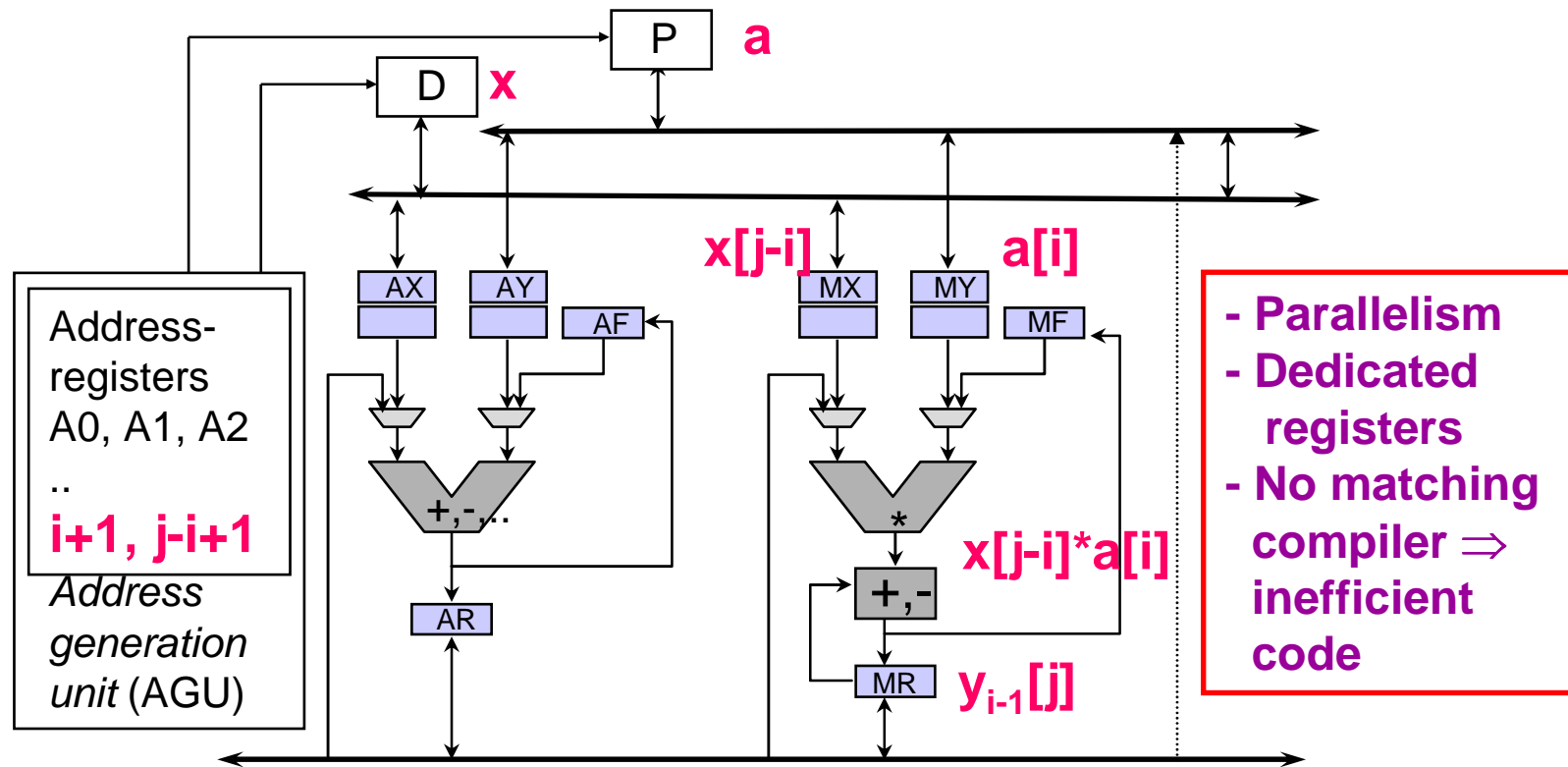Graphics: © Alexandra Nolte, Gesine Marwedel, 2003

# Reason for compiler-problems:
# Application-oriented Architectures

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

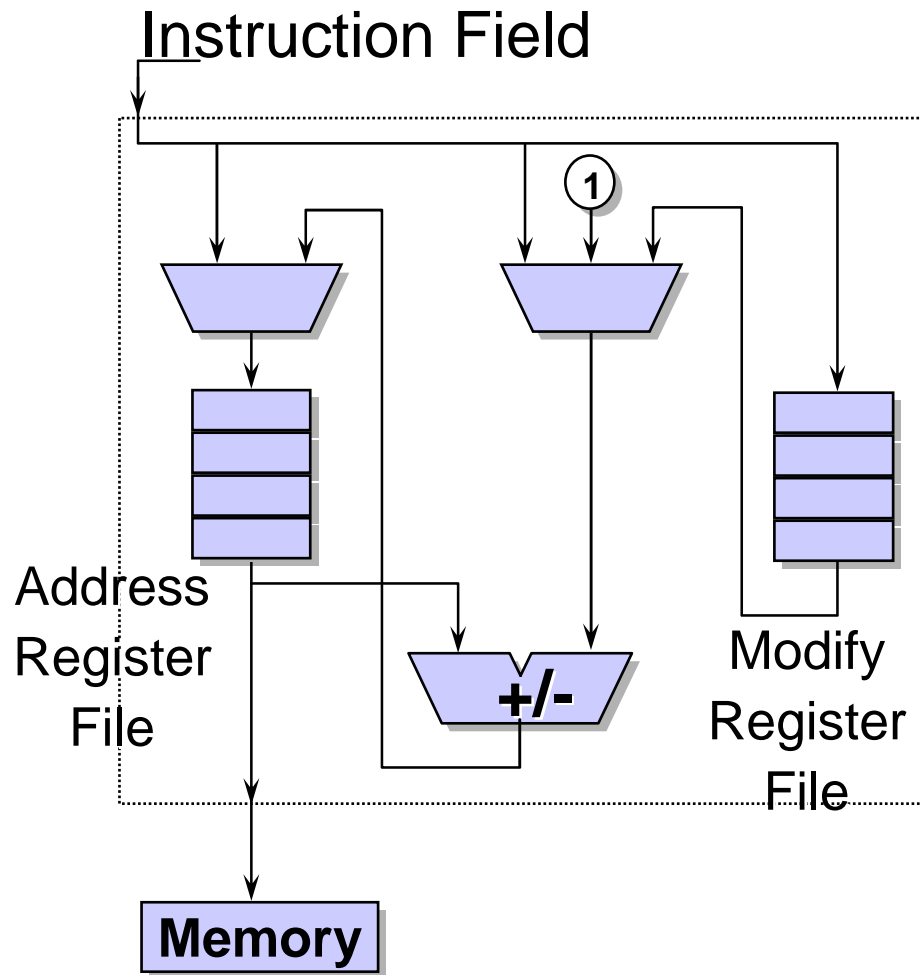$\forall i: 0 \leq i \leq n: y_i[j] = y_{i-1}[j] + x[j-i]*a[i]$

**Architecture:** **Example: Data path ADSP210x**



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

Address-registers A0, A1, A2 ..

**i+1, j-i+1**

*Address generation unit* (AGU)

# Exploitation of parallel address computations

## Generic address generation unit (AGU) model

Instruction Field

Address Register File

+/-

Modify Register File

Memory

**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)**

| | |
|---|---|
| 0 | ar |
| 1 | ai |
| 2 | br |
| 3 | bi |

```
lt ar         ← How to access ar?
mpy br
ltp bi
mpya ar
sacl ar
ltp ai
mpy br
apac
sacl br       ↓ time
```

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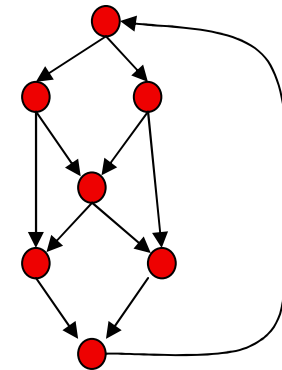
technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 9 -

# 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 :

   $\forall\ v \in V: \sum_{w \in V:(v,w) \in E} g(v,w) = \sum_{w \in V:(w,v) \in E} g(w,v)$   (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.
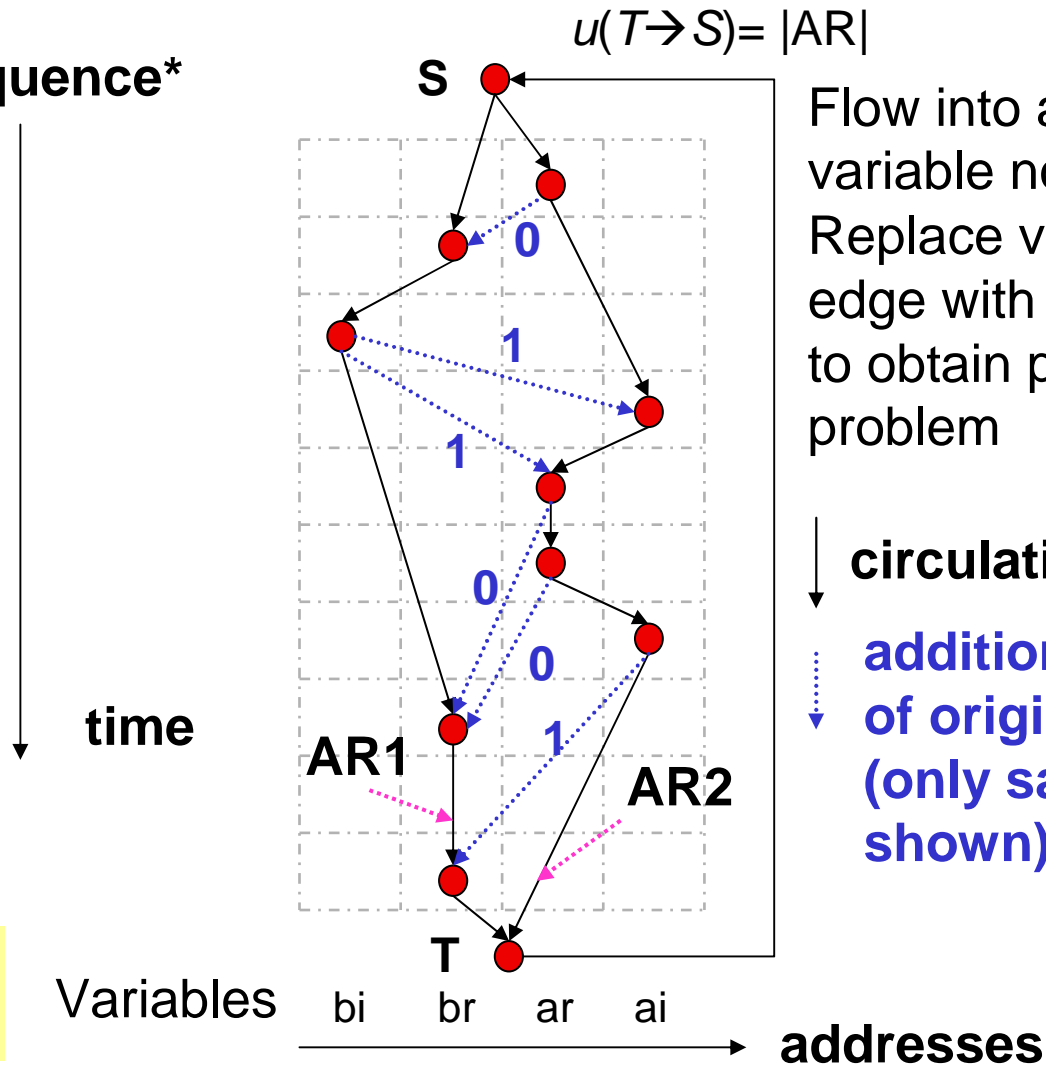
[K.D. Wayne: A Polynomial Combinatorial Algorithm for Generalized Minimum Cost Flow,  http://www.cs.princeton.edu/ ~wayne/ papers/ ratio.pdf]

# Mapping APA to the Minimum Cost Circulation Problem

**Assembly sequence***

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

time

$u(T \rightarrow S) = |AR|$

S

0

1

1

0

0

1

AR1

AR2

T

Variables   bi   br   ar   ai

addresses

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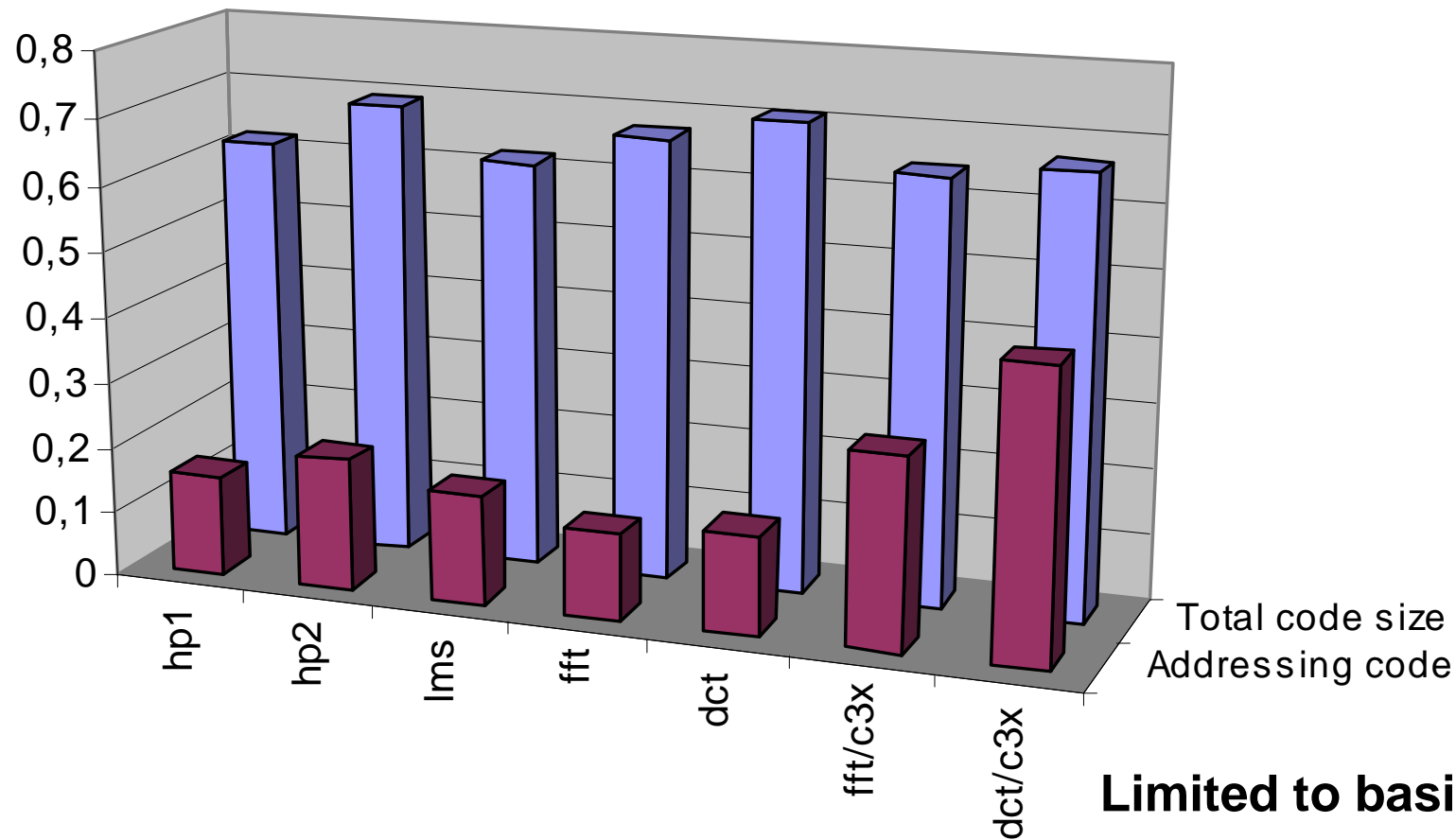
↓ **circulation selected**

**additional edges of original graph (only samples shown)**

\* C2x pro-cessor from ti

# Results according to Gebotys

$$\frac{\text{Optimized code size}}{\text{Original code size}}$$



**Limited to basic blocks**

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 12 -

# Beyond basic blocks:
## - handling array references in loops -
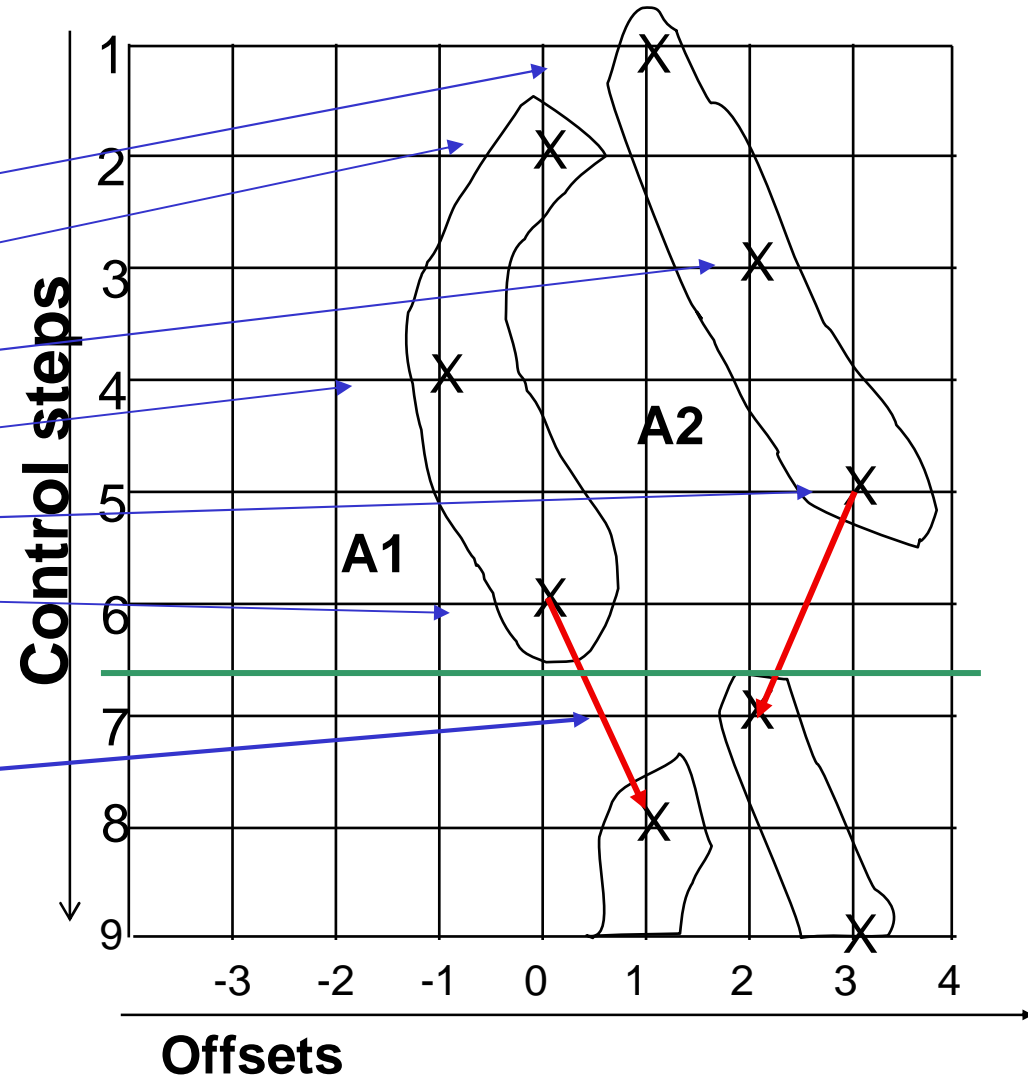
Example:
```
for (i=2; i<=N; i++)
{ .. B[i+1]    /*A2++  */
  .. B[i]      /*A1--  */
  .. B[i+2]    /*A2++  */
  .. B[i-1]    /*A1++  */
  .. B[i+3]    /*A2--  */
  .. B[i]  }   /*A1++  */
```

**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



technische universität dortmund

fakultät für informatik

© p. marwedel, informatik 12, 2011

# Offset assignment problem (OA)
## - Effect of optimised memory layout -

Let's assume that we can modify the memory layout

&#9758; offset assignment problem.

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

&#9758; $k$: number of address registers

&#9758; $m$: number of modify registers

&#9758; $r$: 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).

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12, 2011

- 14 -

# ☞ 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 |
|---|---|
| 1 | b |
| 2 | c |
| 3 | d |

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

cost: 4

| 0 | b |
|---|---|
| 1 | d |
| 2 | c |
| 3 | a |

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

cost: 2

technische universität
dortmund

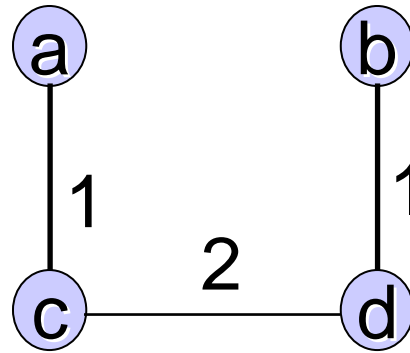fakultät für
informatik

© p. marwedel,
informatik 12, 2011

- 15 -

# SOA example: Access sequence, access graph and Hamiltonian paths

access sequence: b d a c d c

access graph

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

memory layout

| | |
|---|---|
| 0 | b |
| 1 | d |
| 2 | c |
| 3 | a |

SOA used as a building block for more complex situations

➡ significant interest in good SOA algorithms

[Bartley, 1992; Liao, 1995]

technische universität
dortmund

fakultät für
informatik
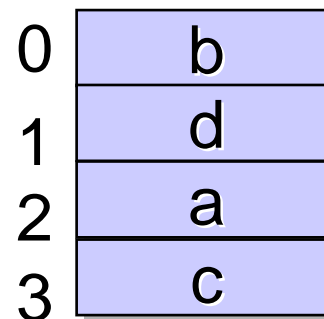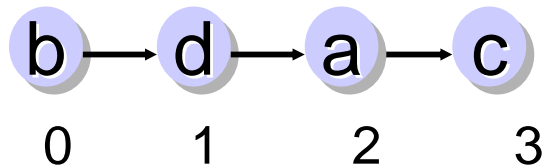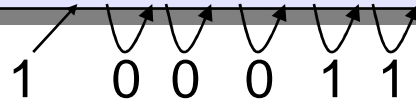
© p. marwedel,
informatik 12,  2011

-  16  -

# Naïve SOA

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

Example:

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

1   0   0   0   1   1

b — d — a — c
0   1   2   3

0  b
1  d
2  a
3  c

memory layout

technische universität
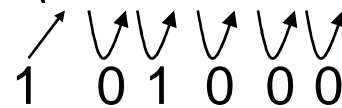dortmund

fakultät für
informatik
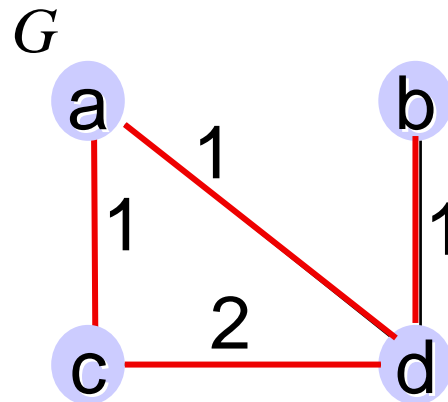
© p. marwedel,
informatik 12, 2011

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

1  0 1 0  0 0

*G*

a     b

1

1       1

Implicit edges of weight 0 for all unconnected nodes.

    2

c     d

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

*G'*

a     b

1       1

    2

c     d

# Liao's algorithm
# on a more complex graph

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

$G$



$G'$

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 19 -

# Further optimizations

Peter Marwedel
TU Dortmund
Informatik 12
Germany

2010/01/13

# Multiple memory banks
## - Sample hardware -



X-Mem

X0
X1
Y0
Y1

Y-Mem

From AGU
or immediate

From AGU
or immediate

Multiplier

ALU

Parallel moves
possible if using
different memory
banks

Shifter

A
B

Shifter

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12, 2011

- 21 -

# Multiple memory banks
# - Constraint graph generation -

**Constraint graph**
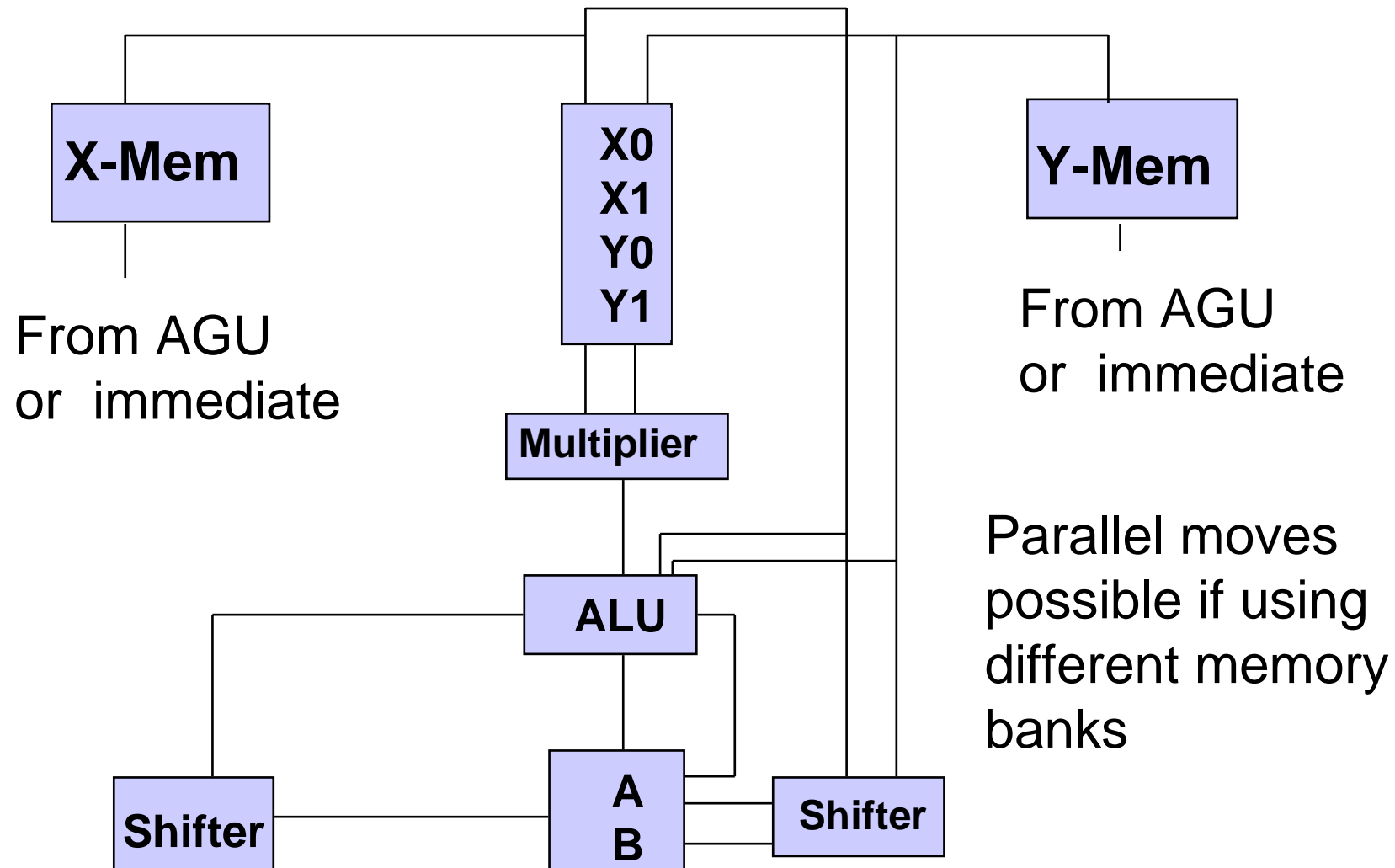
Precompacted
code
(symbolic variables
and registers)

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

v0 — v1 ← ········ {X-Mem, Y-Mem}

r0 ▬▬ r1 ← ········ {X0,X1,Y0, Y1, A, B}

Do not assign to
same register

v2 — v3 ← ········ {X-Mem, Y-Mem}

r2 ▬▬ r3 ← ········ {X0,X1,Y0, Y1, A,B}

Links maintained, more constraints ...

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 22 -

# Multiple memory banks
## Code size reduction through simulated annealing



[Sudarsanam, Malik, 1995]

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12, 2011

- 23 -

# Exploitation of instruction level parallelism (ILP)

Several transfers in the same cycle:

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 24 -

# 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. ~(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)

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

ACCU := u(n - 1)
TR      := e(n - 1)
PR       := TR × K1
TR       := e(n)
e(n - 1) := e(n)
ACCU := ACCU + PR
PR       := TR × K0
ACCU := ACCU + PR
u(n)    := ACCU

ACCU:= u(n - 1)
TR      := e(n - 1)
PR      := TR × K1
e(n - 1):= e(n) || TR:= e(n)  ||
             ACCU:= ACCU + PR
PR      := TR × K0
ACCU:= ACCU + PR
u(n)    := ACCU

- From 9 to 7 cycles through compaction
-

# Exploitation of instruction level parallelism (ILP)

Results obtained through integer programming:

Code size reduction [%]



Compaction times: 2 .. 35 sec

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12, 2011

- 27 -

# 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];
```

**b**

MMAdd (4 x 8/16 bit)

**c**

**a**

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 28 -

# Improvements for M3 DSP due to vectorization

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

# Scheduling for partitioned data paths

Schedule depends on which data path is used.

Cyclic dependency of scheduling and assignment.

'C6x:

Data path A

Data path B

register file A

register file B

| L1 | S1 | M1 | D1 | | D2 | M2 | S2 | L2 |

Address bus

Data bus

# 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 :=
   LISTSCHEDULING(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:=LISTSCHEDULING(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;
```

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 31 -

# 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:***

| | | | | |
|---|---|---|---|---|
| **15 (TriMedia) bzw.** | | | | |
| **40 (C62xx) *delay*** | | | | |
| ***slots*** | | | | |

**Condi-tional jump**

**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**

technische universität
dortmund
fakultät für
informatik
© p. marwedel,
informatik 12, 2011
- 33 -

# 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

if 1

if 2

- 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!

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 35 -

# 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

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 36 -

# 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-statements**

| T1 | T2 |
|----|----|
|    |    |

comparison

if

if

| T1 | T2 |
|----|----|
| T3 | T4 |

< >

**top-down**

if    if

# 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%

technische universität dortmund

fakultät für informatik

© p. marwedel, informatik 12, 2011

# Function inlining: advantages and limitations

## Advantage: low calling overhead

Function sq(c:integer)
 return:integer;
begin
 return c*c
end;
....
a=sq(b);
....

*branching*

*Inlining*

push PC;
push b;
BRA sq;
 pull R1;
 mul R1,R1,R1;
 pull R2;
 push R1;
 BRA (R2)+1;
pull R1;
ST R1,a;

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

## 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



technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 40 -

# 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 | 0000000000000000000000000 |
| 105 | 0010000000110000111011111 |
| 110 | 1011100101110000111111111 |
| 115 | 1011000000000100100011001 |
| 120 | 1011010010100010011011110 1 |
| 125 | 1011000000101000010011110 1 |
| 130 | 0011000000001010010011100 0 |
| 135 | 1011001000111010111011110 1 |
| 140 | 1011101111111010111111111 |
| 145 | 1011011010101010011011110 1 |
| 150 | 1011011000001011011011110 1 |

**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

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 43 -

# Retargetable Compilers vs. Standard Compilers



Standard Compiler: Machine model → mainly manual process → Compiler

Retargetable Compiler: Machine model → mainly automatic process → Compiler

**Developer retargetability:** compiler specialists responsible for retargeting compilers.

**User retargetability:** users responsible for retargeting compiler.

# Compiler structure

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 45 -

# Code selection = covering DFGs



load     load     load     load

MEM    MEM    MEM    MEM

\*     mul     \*    mul

mac     +    add    mac

Does not yet consider data moves to input registers.

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 46 -

# Code selection by tree parsing (1)

Specification of grammar for generating a iburg parser*:

terminals: {MEM, *, +}

non-terminals: {reg1,reg2,reg3}

start symbol: reg1

rules:

"add" (cost=2): reg1 -> + (reg1, reg2)

"mul" (cost=2): reg1 -> * ( reg1,reg2)

"mac" (cost=3): reg1 -> + (*(reg1,reg2), reg3)

"load" (cost=1): reg1 -> MEM

"mov2"(cost=1): reg2 -> reg1

"mov3"(cost=1): reg3 -> reg1

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12, 2011

- 47 -

# Code selection by tree parsing (2)
## - nodes annotated with (register/pattern/cost)-triples -

**"load"(cost=1):**
    **reg1 -> MEM**
**"mov2"(cost=1):**
    **reg2 -> reg1**
**"mov3"(cost=1):**
    **reg3 -> reg1**

**"add" (cost=2):**
    **reg1 -> +(reg1, reg2)**
**"mul" (cost=2):**
    **reg1 -> *(reg1,reg2)**
**"mac" (cost=3):**
    **reg1->+(*(reg1,reg2), reg3)**

reg1:load:1
reg2:mov2:2
reg3:mov3:2

MEM   MEM   MEM   MEM

reg1:mul:5
reg2:mov2:6
reg3:mov3:6

*        *

+

reg1:add:13
reg1:mac:12

# Code selection by tree parsing (3)
## - final selection of cheapest set of instructions -

load        load        load        load

MEM         MEM         MEM         MEM

mov2

mov2

*           *

mul

Includes
routing of
values between     mac
various
registers!              +

mov3

reg1:add:13
reg1:mac:12 ⟵⟵⟵

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12, 2011

- 49 -

# Dynamic Voltage Scaling

Peter Marwedel
TU Dortmund
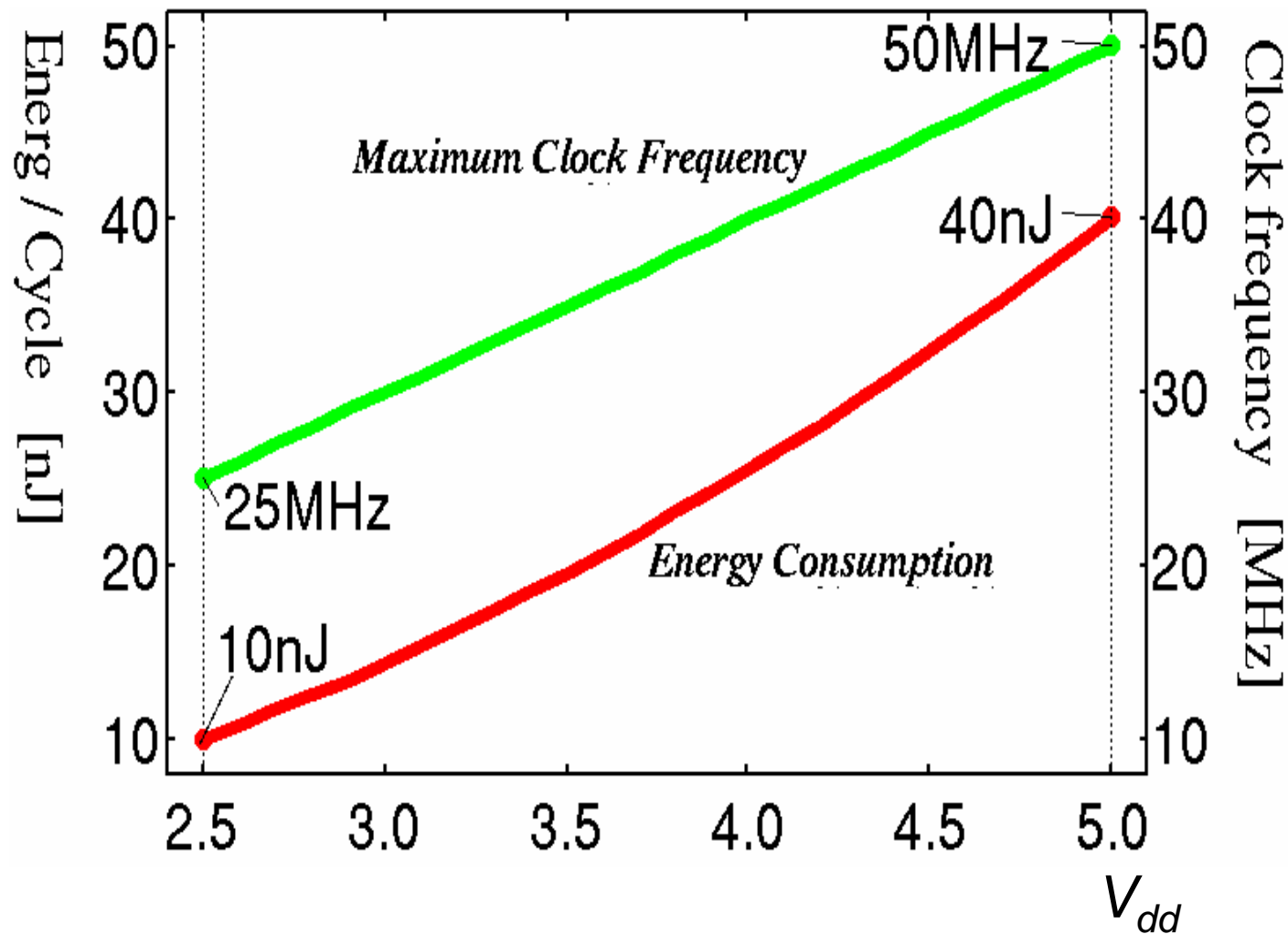Informatik 12
Germany

2009/01/17

# Voltage Scaling and Power Management
## Dynamic Voltage Scaling

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 51 -

# Recap from chapter 3: Fundamentals of dynamic voltage scaling (DVS)

Power consumption of CMOS circuits (ignoring leakage):

$$P = \alpha \; C_L \, V_{dd}^2 \; f \;\; \text{with}$$

$\alpha$ : switching activity

$C_L$ : load capacitance

$V_{dd}$ : supply voltage

$f$ : clock frequency

Delay for CMOS circuits:

$$\tau = k \, C_L \frac{V_{dd}}{\left(V_{dd} - V_t\right)^2} \;\; \text{with}$$

$V_t$ : threshhold voltage

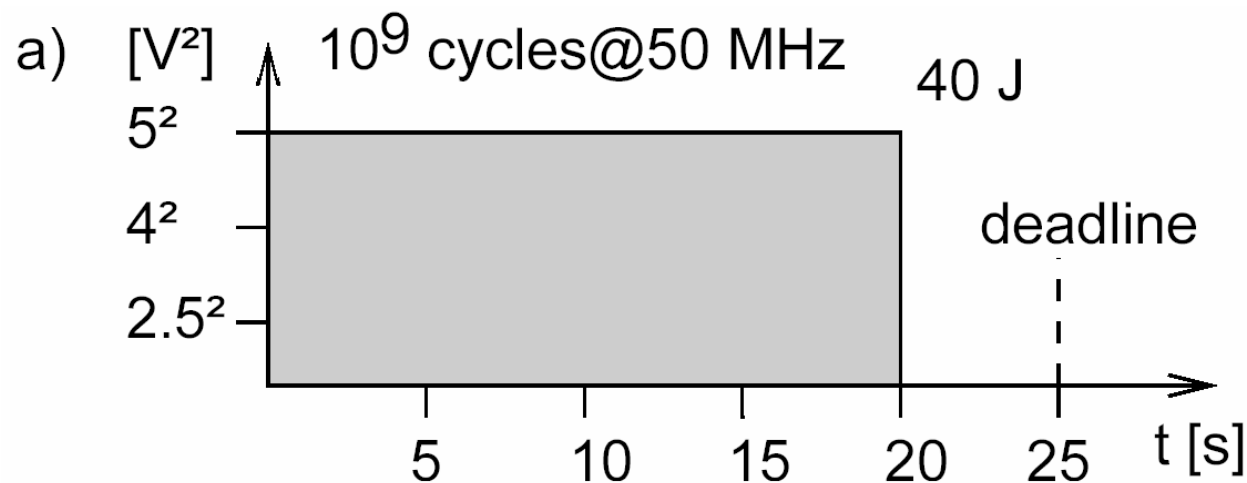$(V_t \text{ substancially} < \text{than } V_{dd})$

☞ Decreasing $V_{dd}$ reduces $P$ quadratically, while the run-time of algorithms is only linearly increased (ignoring the effects of the memory system).

# Example: Processor with 3 voltages
# Case a): Complete task ASAP

Task that needs to execute $10^9$ cycles within 25 seconds.

| $V_{dd}$ [V] | 5.0 | 4.0 | 2.5 |
|---|---|---|---|
| Energy per cycle [nJ] | 40 | 25 | 10 |
| $f_{max}$ [MHz] | 50 | 40 | 25 |
| cycle time [ns] | 20 | 25 | 40 |

a) [V²] $10^9$ cycles@50 MHz 40 J

$E_a = 10^9 \times 40 \times 10^{-9}$ [J]
  = 40 [J]

# Case b): Two voltages

| $V_{dd}$ [V] | 5.0 | 4.0 | 2.5 |
|---|---|---|---|
| Energy per cycle [nJ] | 40 | 25 | 10 |
| $f_{max}$ [MHz] | 50 | 40 | 25 |
| cycle time [ns] | 20 | 25 | 40 |

b) [V²] 750M cycles @ 50 MHz + 250M cycles @ 25

32.5 J    deadline

$E_b$= 750 10⁶ x 40 10 ⁻⁹ +
250 10⁶ x 10 10⁻⁹ [J]

= 32.5 [J]

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12, 2011

- 54 -

# Case c): Optimal voltage

| $V_{dd}$ [V] | 5.0 | 4.0 | 2.5 |
|---|---|---|---|
| Energy per cycle [nJ] | 40 | 25 | 10 |
| $f_{max}$ [MHz] | 50 | 40 | 25 |
| cycle time [ns] | 20 | 25 | 40 |

c) [V²]

$10^9$ cycles@40 MHz

25 J

$5^2$

$4^2$

$2.5^2$

5    10    15    20    25    t [s]

$E_c = 10^9 \times 25 \times 10^{-9}$ [J]

$= 25$ [J]

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12, 2011

- 55 -

# Observations

- A minimum energy consumption is achieved for the ideal supply voltage of 4 Volts.

- In the following: **variable voltage processor =** processor that allows **any** supply voltage up to a certain maximum.

- It is expensive to support truly variable voltages, and therefore, actual processors support only a few fixed voltages.

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 56 -

# Generalisation

Lemma [Ishihara, Yasuura]:

- If a variable voltage processor completes a task before the deadline, the energy consumption can be reduced.
- If a processor uses a single supply voltage $V$ and completes a task $T$ just at its deadline, then $V$ is the unique supply voltage which minimizes the energy consumption of $T$.
- If a processor can only use a number of discrete voltage levels, then the two voltages which (almost*) minimize the energy consumption are the two immediate neighbors of the ideal voltage $V_{ideal}$ possible for a variable voltage processor.

  ----------------

  * Except for small amounts of energy resulting from the fact that cycle counts must be integers.

# The case of multiple tasks:
# Assignment of optimum voltages to a set of tasks

$N$ : the number of tasks

$EC_j$ : the number of executed cycles of task $j$

$L$ : the number of voltages of the target processor

$V_i$ : the $i^{th}$ voltage, with $1 \leq i \leq L$

$F_i$ : the clock frequency for supply voltage $V_i$

$T$ : the global deadline at which all tasks must have been completed

$X_{i,j}$ : the number of clock cycles task $j$ is executed at voltage $V_i$

$SC_j$ : the average switching capacitance during the execution of task $j$ ($SC_j$ comprises the actual capacitance $CL$ and the switching activity $\alpha$)

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12, 2011

- 58 -

# Designing an ILP model

Simplifying assumptions of the ILP-model include the following:

- There is one target processor that can be operated at a limited number of discrete voltages.

- The time for voltage and frequency switches is negligible.

- The worst case number of cycles for each task are known.

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 59 -

# Energy Minimization
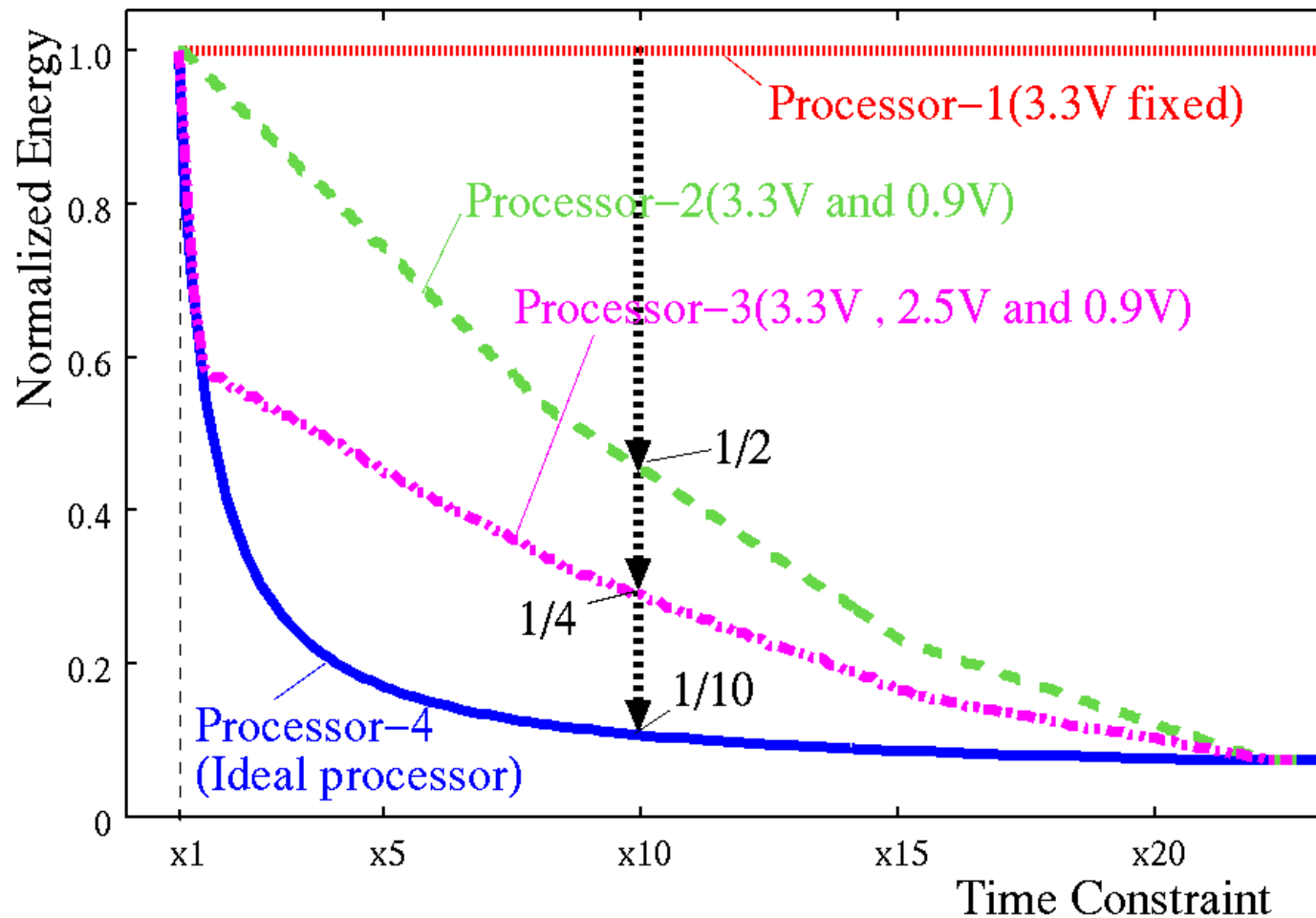# using an Integer Linear Programming Model

Minimize

$$E = \sum_{j=1}^{N} \sum_{i=1}^{L} SC_j \cdot x_{i,j} \cdot V_i^2$$

subject to

$$\forall j : \sum_{i=1}^{L} x_{i,j} = EC_j$$

and

$$\sum_{i=1}^{L} \sum_{j=1}^{N} \frac{x_{i,j}}{F_i} \leq T$$

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 60 -

# Experimental Results

3 tasks; $EC_j$=50 $10^9$; $SC_1$=50 pF; $SC_2$=100 pF; $SC_3$=150 pF

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12, 2011

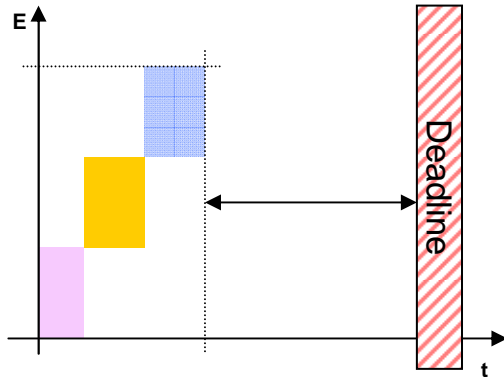© Yasuura, 2000

- 61 -

# Task-level concurrency management

- The dynamic behavior of applications getting more attention.

- Energy consumption reduction is the main target.

- Some classes of applications (i.e. video processing) have a considerable variation in processing power requirements depending on input data.

- Static design-time methods becoming insufficient.

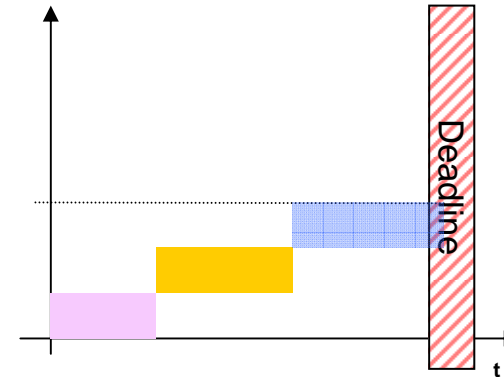- Runtime-only methods not feasible for embedded systems.

→ How about mixed approaches?

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12, 2011

- 62 -

# Example of a mixed TCM
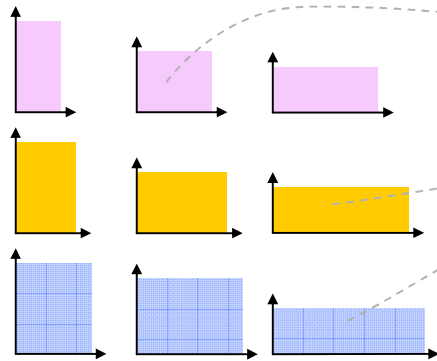
Task1

Task2

Task3
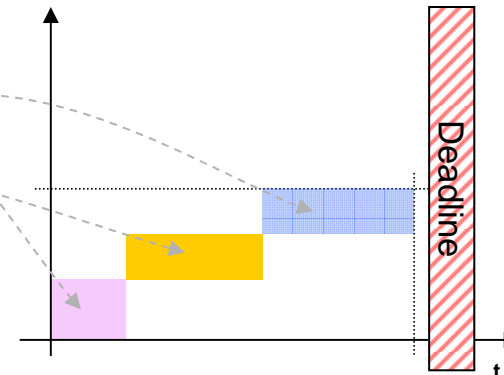


Static (compile-time) methods can ensure WCET feasible schedules, but waste energy in the average case.

…or they can define a probability for violating the deadline.

Mixed methods use compile-time analysis to define a set of possible execution parameters for each task.

Runtime scheduler selects the most energy saving, deadline preserving combination.

**[IMEC, Belgium, http://www.imec.be/]**

technische universität dortmund

fakultät für informatik

# Dynamic power management (DPM)

Dynamic Power management tries to assign optimal power saving states.

- Questions: When to go to an power-saving state? Different, but typically complex models:

- Markov chains, renewal theory , ….

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 64 -

# Summary

- Worst-case execution time aware compilation

- The offset assignment problem

- Retargetability and code selection

- Dynamic voltage scaling (DVS)

  - An ILP model for voltage assignment in a multi-tasking system

- Dynamic power management (DPM) (briefly)

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12,  2011

- 65 -