

UNIVERSITÄT DORTMUND



# Memory-architecture aware compilation

Peter Marwedel

University of Dortmund  
Informatik 12  
44221 Dortmund, Germany

[peter.marwedel@udo.edu](mailto:peter.marwedel@udo.edu)

<http://ls12-www.cs.uni-dortmund.de>

phone: +49 (231) 755 6111

mobile: +49 (173) 5165815

fax: +49 (231) 755 6116

UNIVERSITÄT DORTMUND

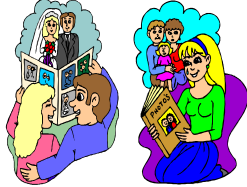


# Introduction


Universität Dortmund

## The Problem with Memories



Memories?



Oops!  
Memories!



1. Increasing speed gap
2. Major consumer of electrical energy
3. Timing predictability difficult to achieve
4. ...

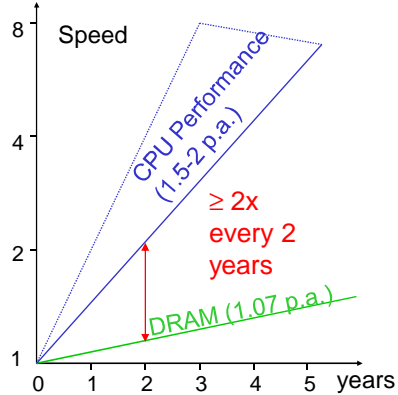
Or: Why work on processors if memory is where the bottleneck is?

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 3 - H/PEAC artit

Universität Dortmund


## Trends for the speeds

Speed gap between processor and main DRAM increases



Similar problems also for embedded systems & MPSoCs

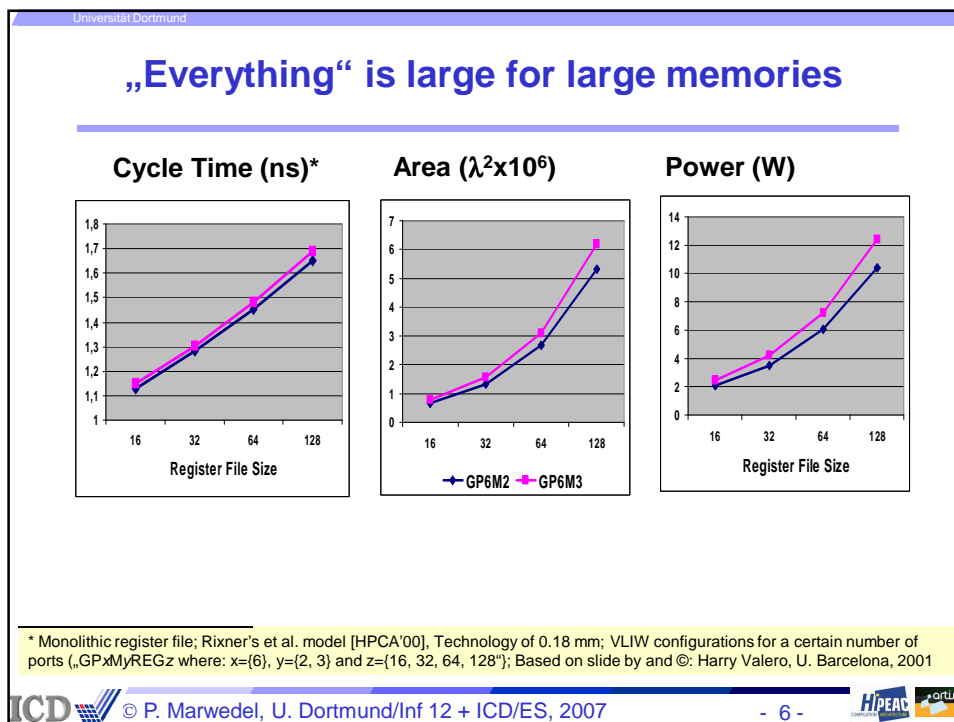
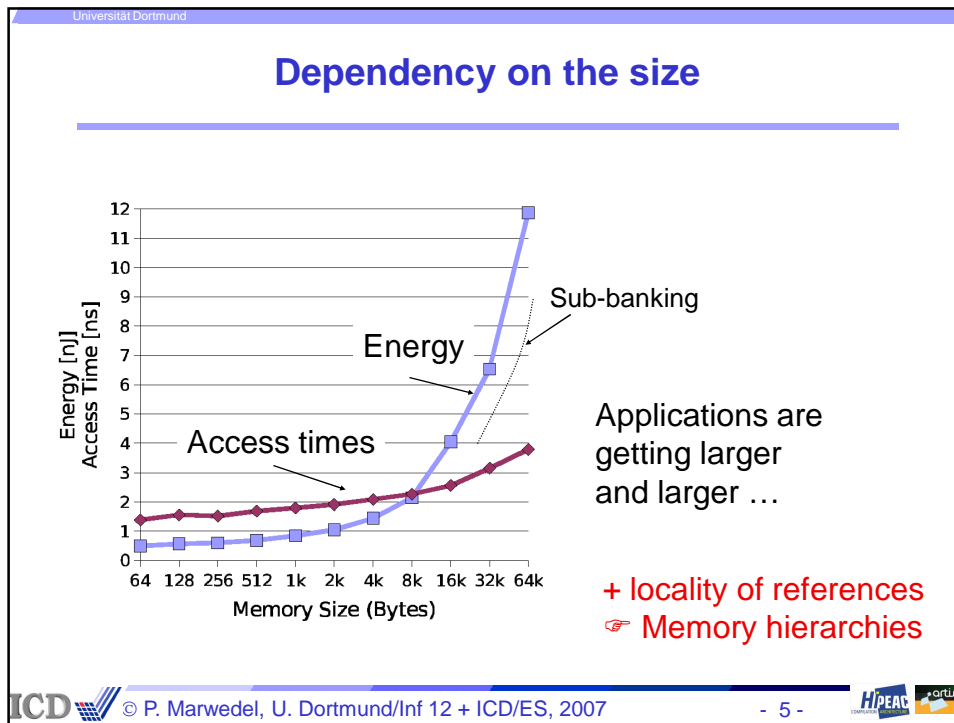
In the future:  
Memory access times >> processor cycle times  
"Memory wall" problem

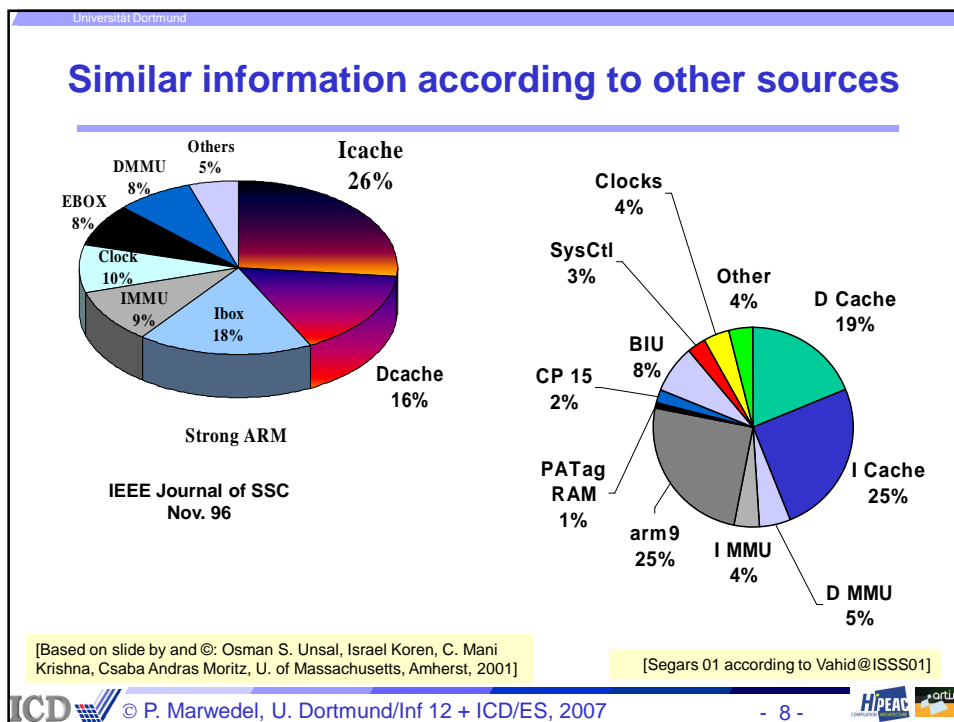
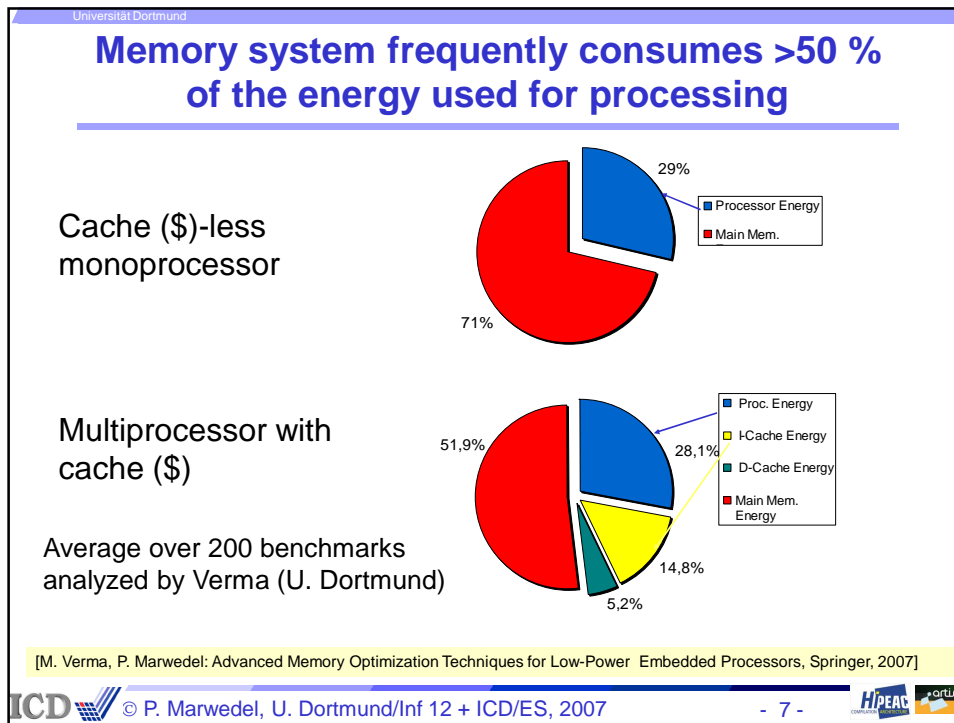


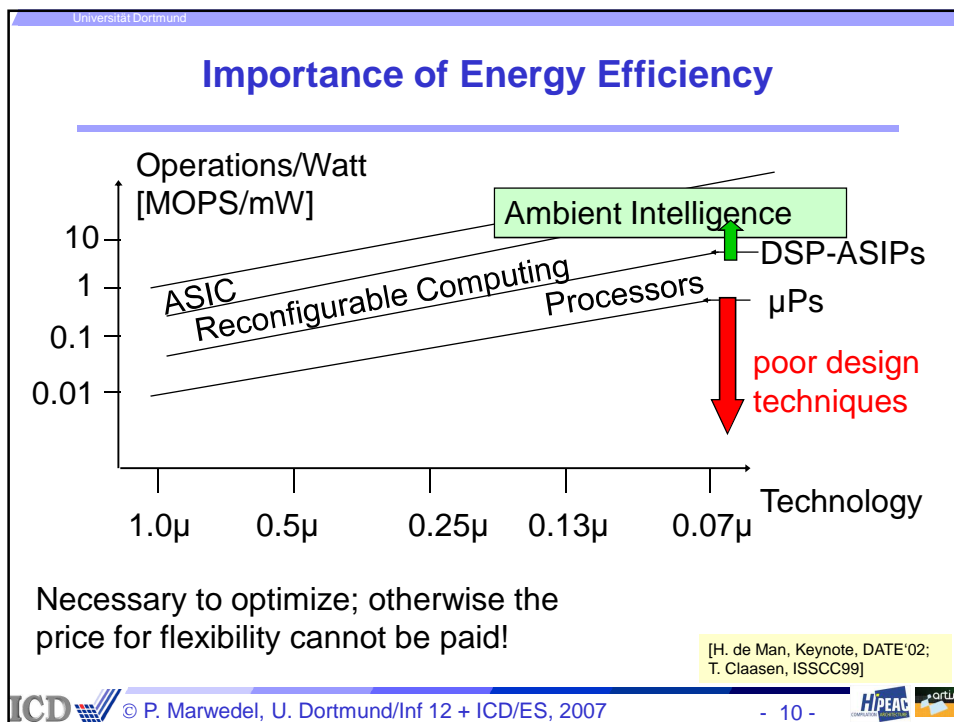
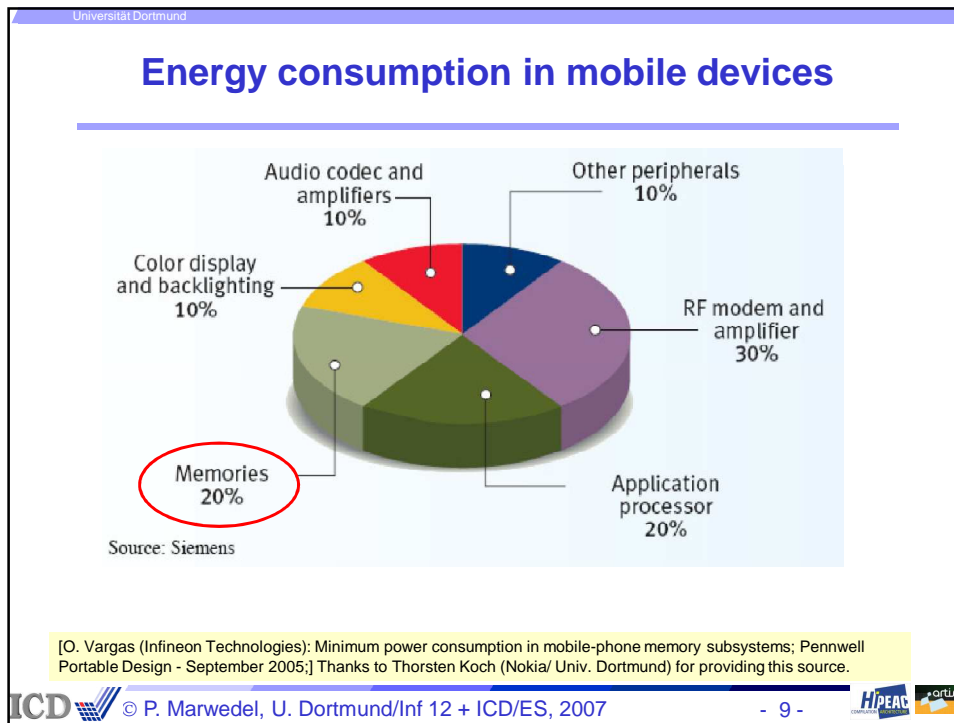
References integrated into slides for modularity & context

[P. Machanik: Approaches to Addressing the Memory Wall, TR Nov. 2002, U. Brisbane]

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 4 - H/PEAC artit







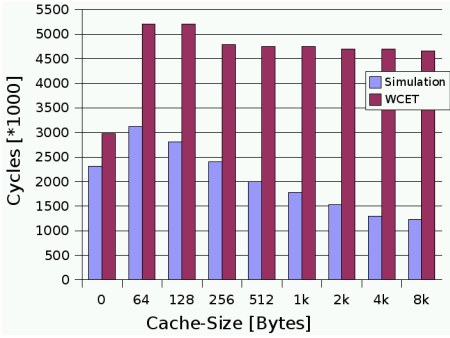
Universitat Dortmund

## Timing Predictability

Many embedded systems are real-time systems

- ☞ computations to be finished in a given amount of time

Most memory hierarchies (e.g. caches) for PC-like systems designed for good average case, not for good worst case behavior.






Cache-Size [Bytes]	Simulation Cycles [*1000]	WCET Cycles [*1000]
0	~2300	~3000
64	~3100	~5200
128	~2800	~5200
256	~2400	~4800
512	~2000	~4800
1k	~1800	~4800
2k	~1500	~4800
4k	~1300	~4800
8k	~1200	~4800

Worst case execution time (WCET) larger than without cache







G.721: using unified Cache@ARM7TDMI




See later slide for experimental setup

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 11 -  

Universitat Dortmund

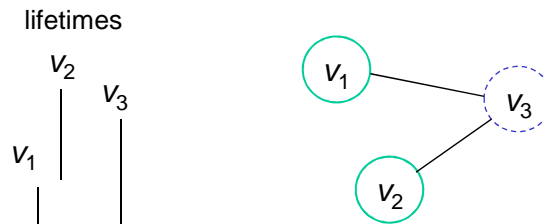
## ☞ Multiple Objectives for Memory System Design

- (Average) Performance
  - Throughput 
  - Latency 
- Energy consumption 
- Predictability, good worst case execution time bound (WCET) 
- Size 
- Cost 
- ....

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 12 -  

## Register allocation

- Registers = fastest level in the memory hierarchy
- ☞ Interest in good global register allocation techniques
- Frequently based on coloring of interference graph



- Registers neither suitable for  $|\text{objects}| \gg \text{single words}$  nor for code
- ☞ Register allocation not considered in the remainder

## Outline

- Introduction, Motivation
- ➡ Scratch-pad memories
  - Characteristics
  - Non-overlying allocation, partitions, stack, heap, WCET
  - Overlaying (dynamic) allocation
  - Multiple processes & processors
- Caches
  - Predictability, loop transformations
  - Prefetching, layout transformations
- Other levels of the memory hierarchy
  - Array folding
  - ADL-driven compilation for hierarchies
  - Rewriting algorithms for memory hierarchies
- Conclusion, Additional references



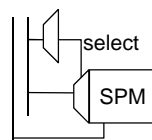
## Scratch pad memories

- Characteristics
- Compilation (1)

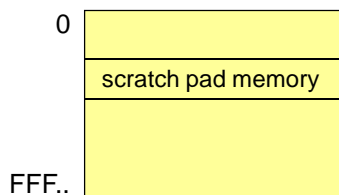
### Scratch pad memories (SPM): Fast, energy-efficient, timing-predictable

SPMs are small,  
physically  
separate  
memories  
mapped into the  
address space;

Selection is by  
an appropriate  
address decoder  
(simple!)



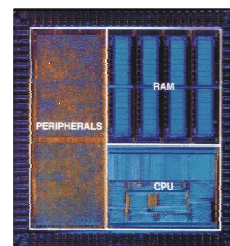
Address space



Small; no  
tag memory

Example

ARM7TDMI  
cores, well-  
known for low  
power  
consumption





## Predictability and scratch-pad memories

... **pre-run-time scheduling** is often the only practical means of providing predictability in a complex system.

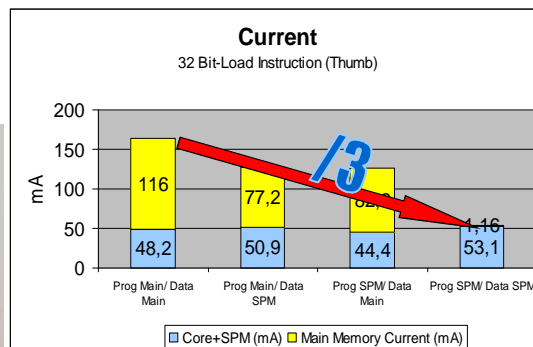
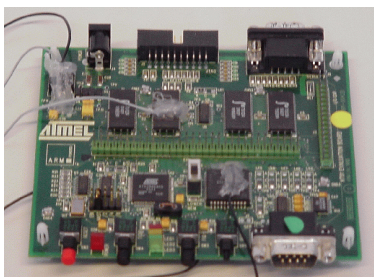
[J. Xu, D. Parnas: On satisfying timing constraints in hard real-time systems, *IEEE Trans. Soft. Engineering*, 1993, p. 70–84]

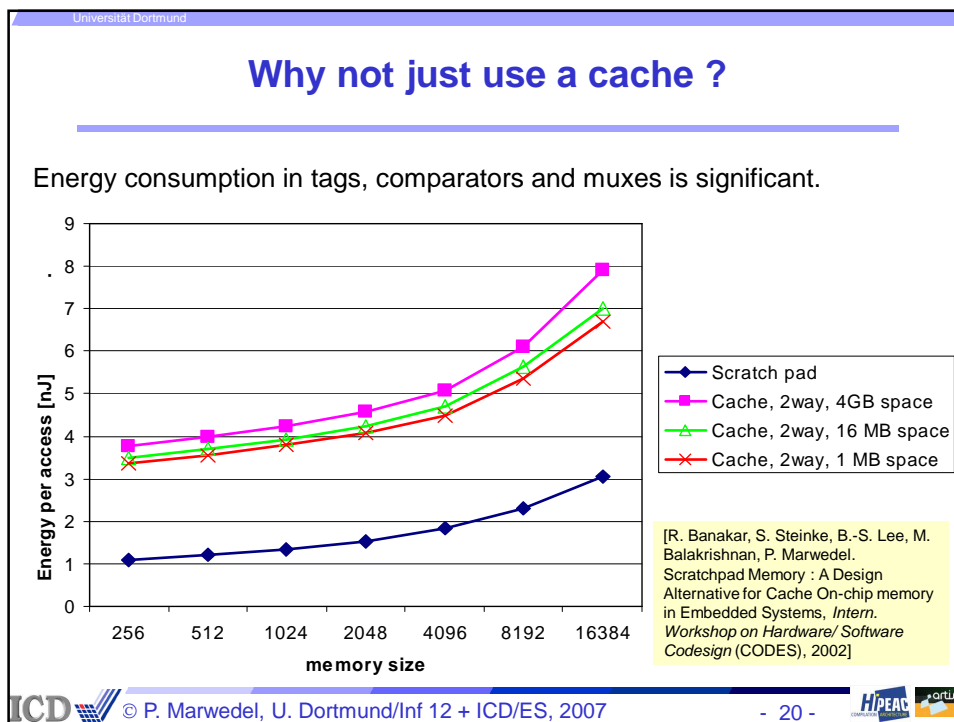
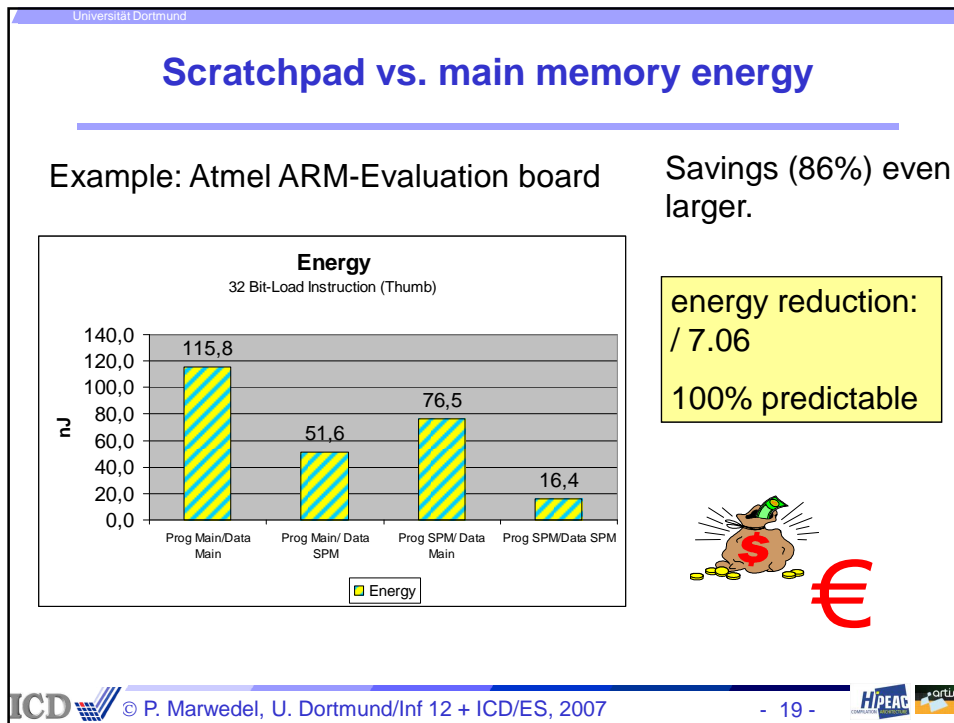
... In essence, we must reinvent computer science. Fortunately, we have quite a bit of knowledge and experience to draw upon. Architecture techniques such as **software-managed caches** promise to deliver much of the benefit of memory hierarchy without the timing unpredictability.

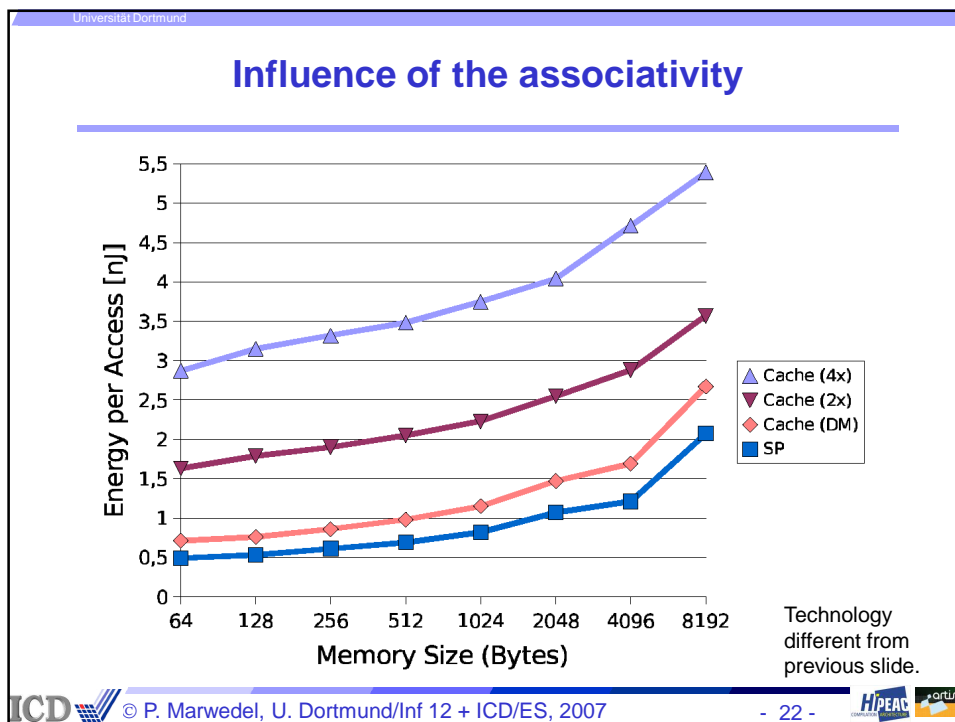
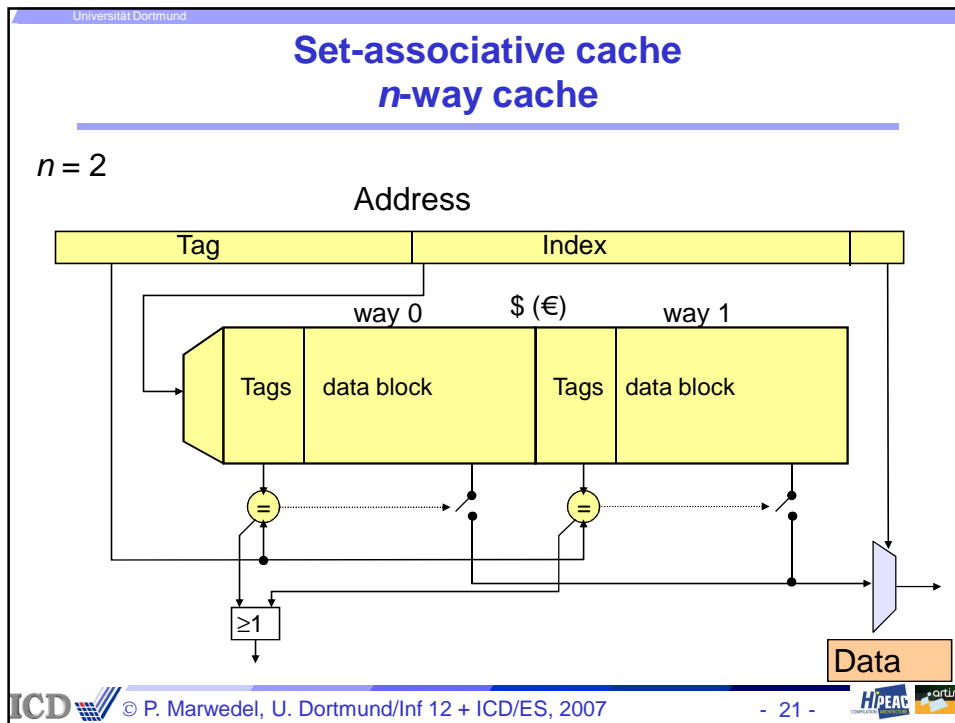
[Ed Lee: Absolutely Positively on Time: What would it take?, *IEEE Computer*, 2005]

## Comparison of currents using measurements

E.g.: ATMEL board with  
ARM7TDMI and  
ext. SRAM










Universität Dortmund

## Availability of SPMs=(“Tightly Coupled Memories”)

ARM CPU Core	Caches	TCM Available
ARM 1026EJ-S	Variable	yes
ARM 1136J(F)-S	Variable	yes
ARM 1176JZ(F)-S	Variable	yes
ARM 926EJ-S	Variable	yes
ARM 1026EJ-S	Variable	yes
ARM 1156T2(F)-S	Variable	yes
ARM 946E-S	Variable	yes
ARM 966E-S	-	yes
ARM 968E-S	-	yes
All others		no

Source: [http://www.arm.com/products/CPUs/core\\_selector.html](http://www.arm.com/products/CPUs/core_selector.html)

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 23 -  

Universität Dortmund

## Current usage for ARM

- Use pragma in C-source to allocate to specific section:**  
For example:  

```
#pragma arm section rwddata = "foo", rodata = "bar"
int x2 = 5; // in foo (data part of region)
int const z2[3] = {1,2,3}; // in bar
```
- Input scatter loading file to linker for allocating section to specific address range**

```
LOAD_ROM_1 0x000
{
  EXEC_ROM_1 0x000
  {
    program1.o (+RO)
  }
  DRAM 0x18000 0x8000
  {
    program1.o (+RW,+ZI)
  }
}
```

Scatter description

Load region description

Execution region description




Input section description

Execution region description

Input section description

(2 different examples)

[http://www.arm.com/documentation/Software\\_Development\\_Tools/index.html](http://www.arm.com/documentation/Software_Development_Tools/index.html)

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 24 -  

## Popular among designers?

Information received from designers:

- Used for buffering speech in mobile phones (mobile phone company)
- Manual mapping of frequently accessed data
- Essentially no idea on how to exploit it (WLAN/Bluetooth specialists & major vendor)

Why not change this situation?

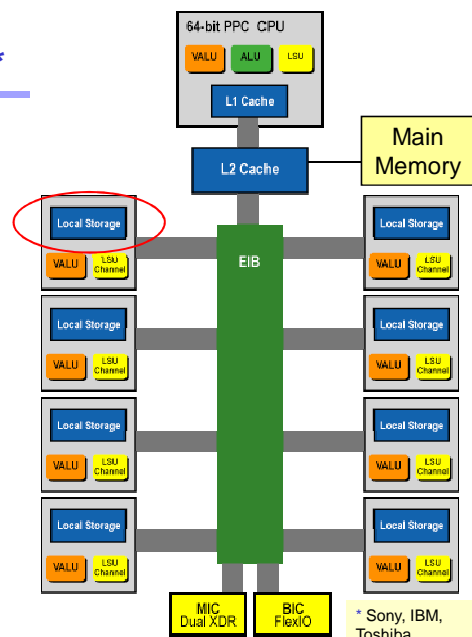
## Similar Concept for the Cell processor \*

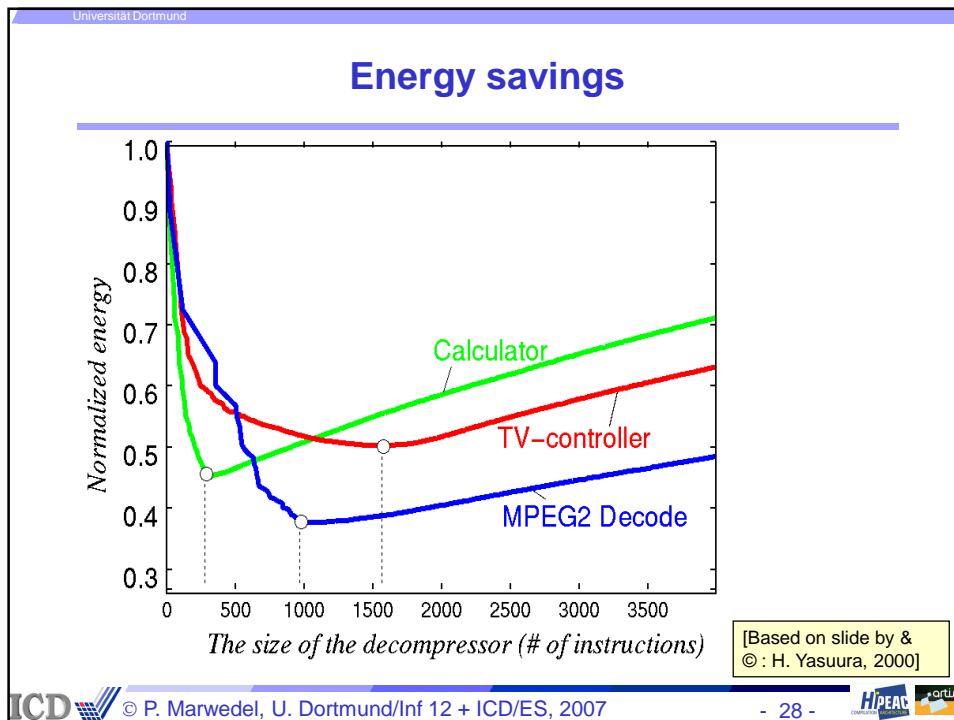
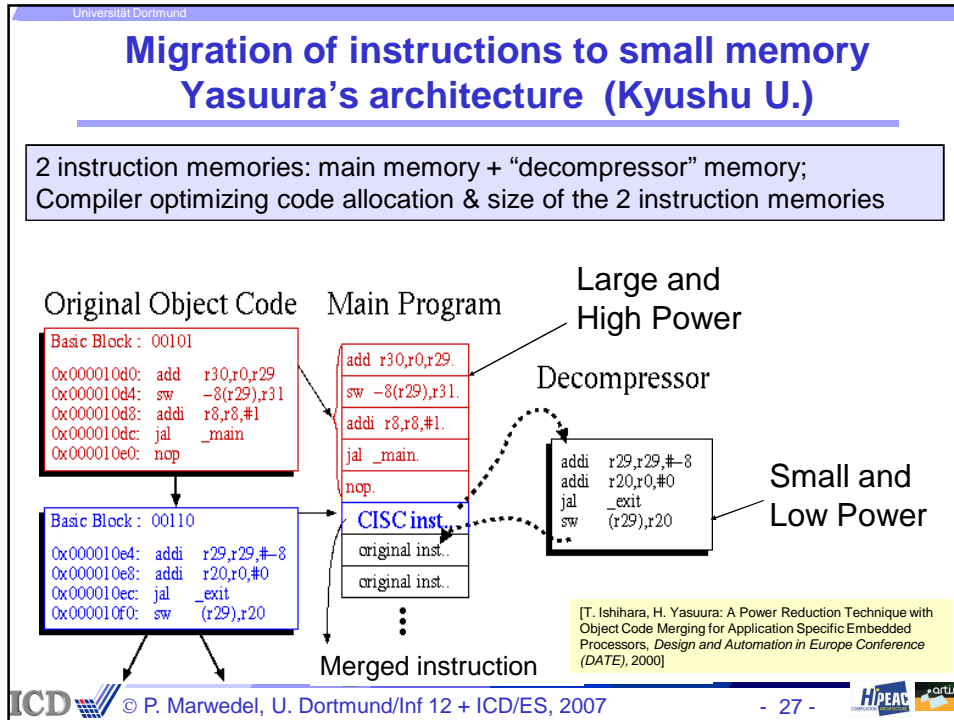
Local SPE processors fetch instructions and data from local storage LS (256 kB). LS **not** designed as a cache. Separate DMA transfers required to fill and spill.

Motivation same as for this tutorial:

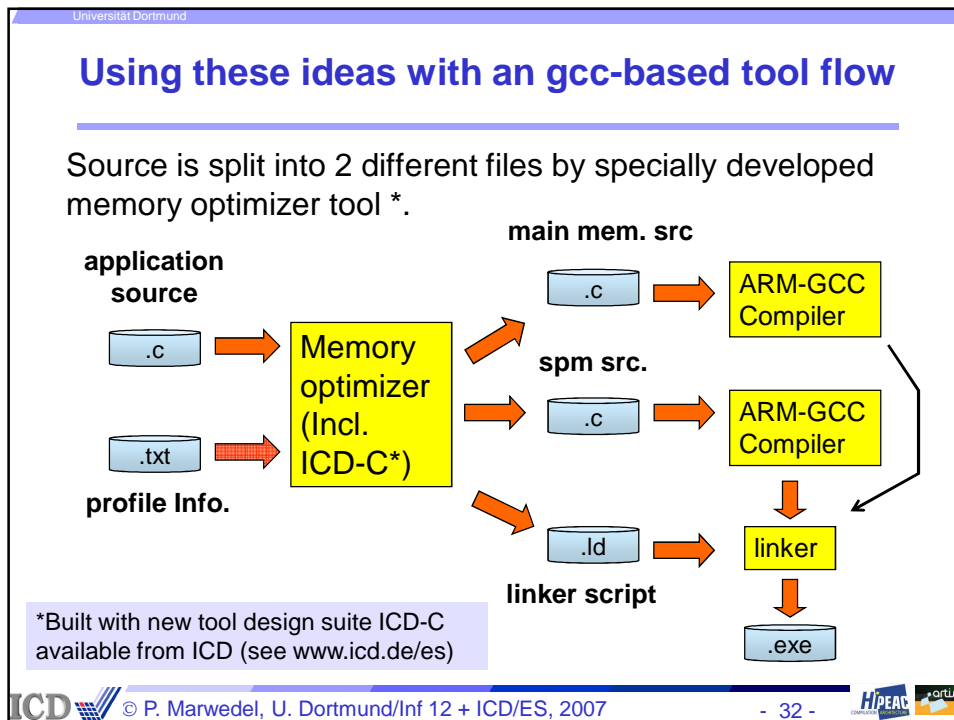
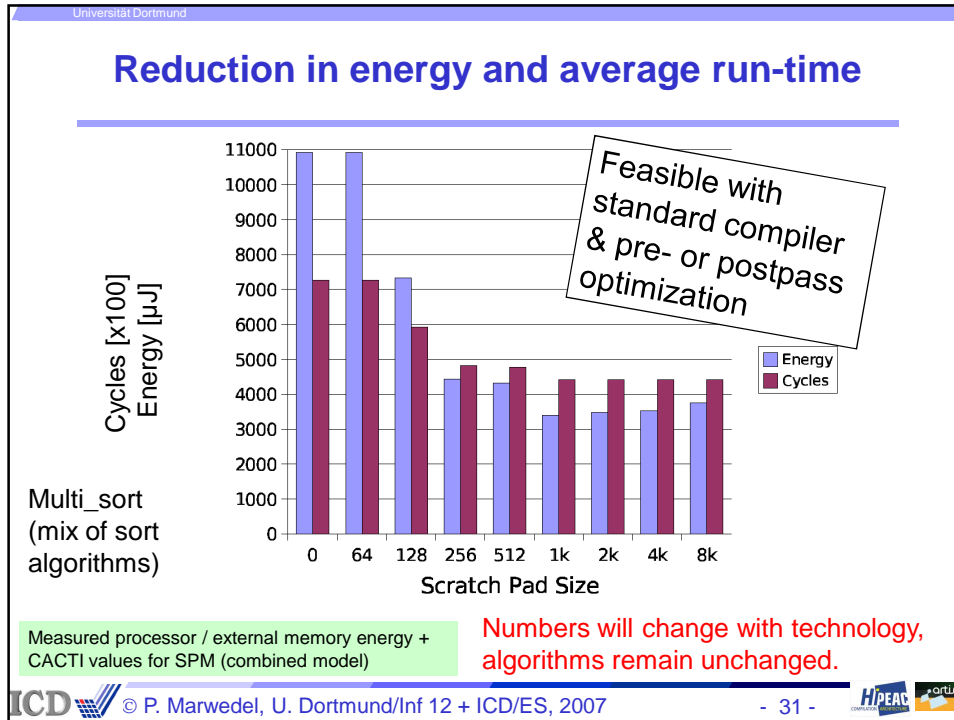
- Large memory latency
- Huge overhead for automatically managed caches

Similar for Infineon TriCore











Universität Dortmund

## Allocation of basic blocks

Fine-grained granularity smoothens dependency on the size of the scratch pad.

Requires additional jump instructions to return to "main" memory.

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 33 -

Universität Dortmund

## Taking consecutive basic blocks into account

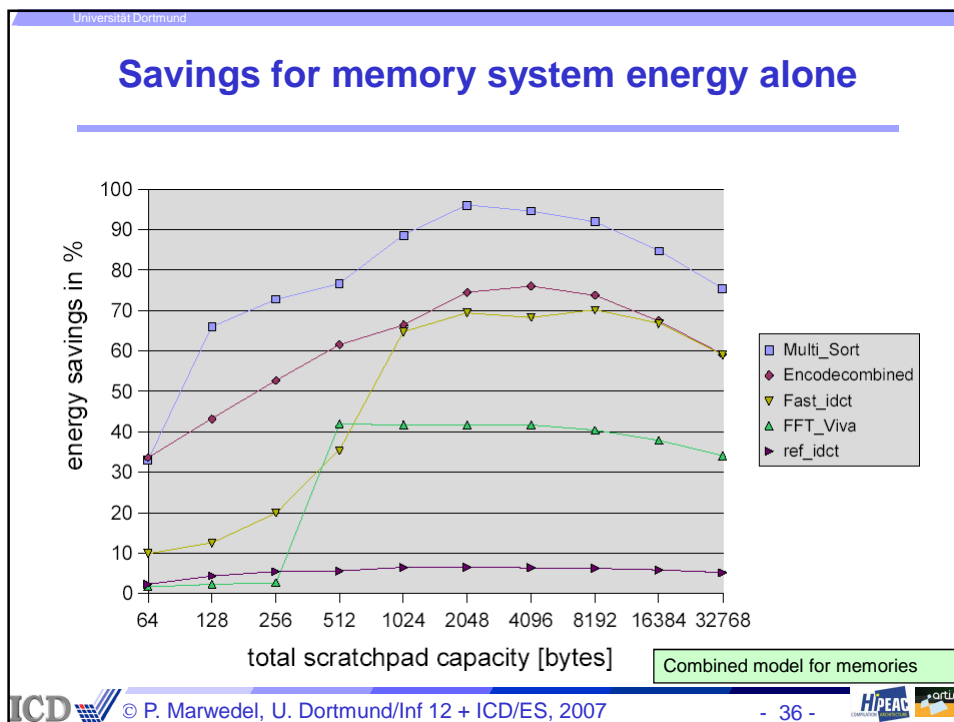
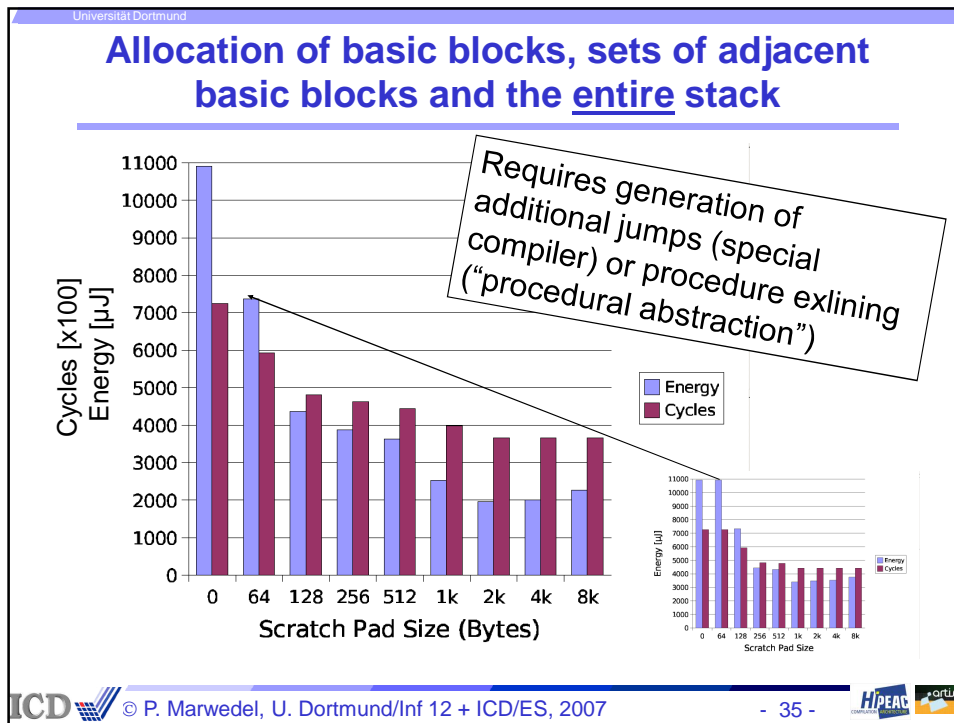
**Approach:**

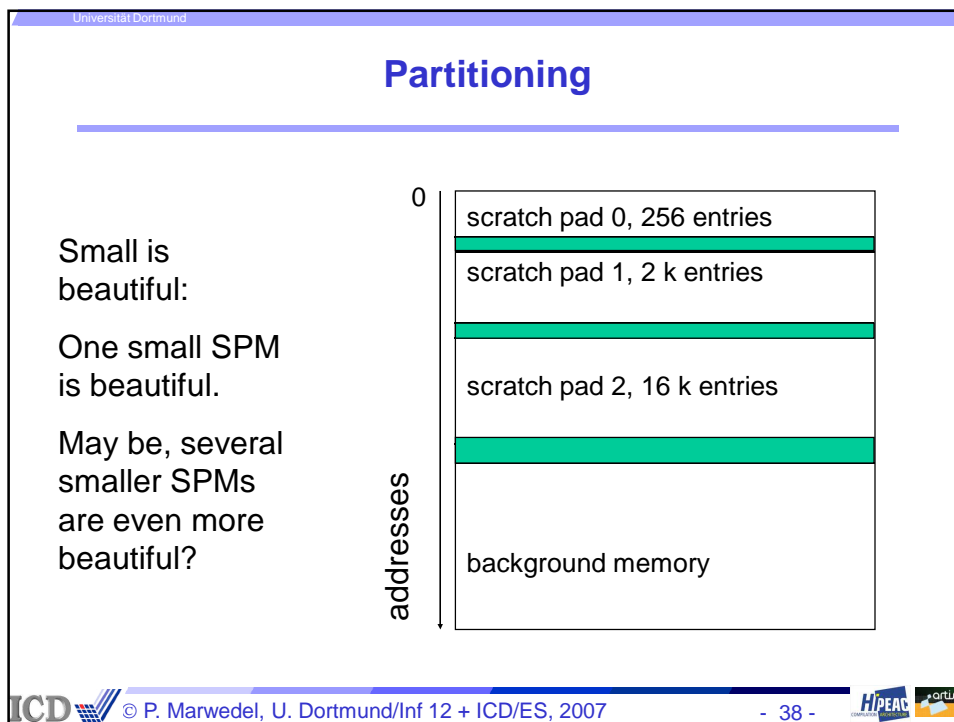
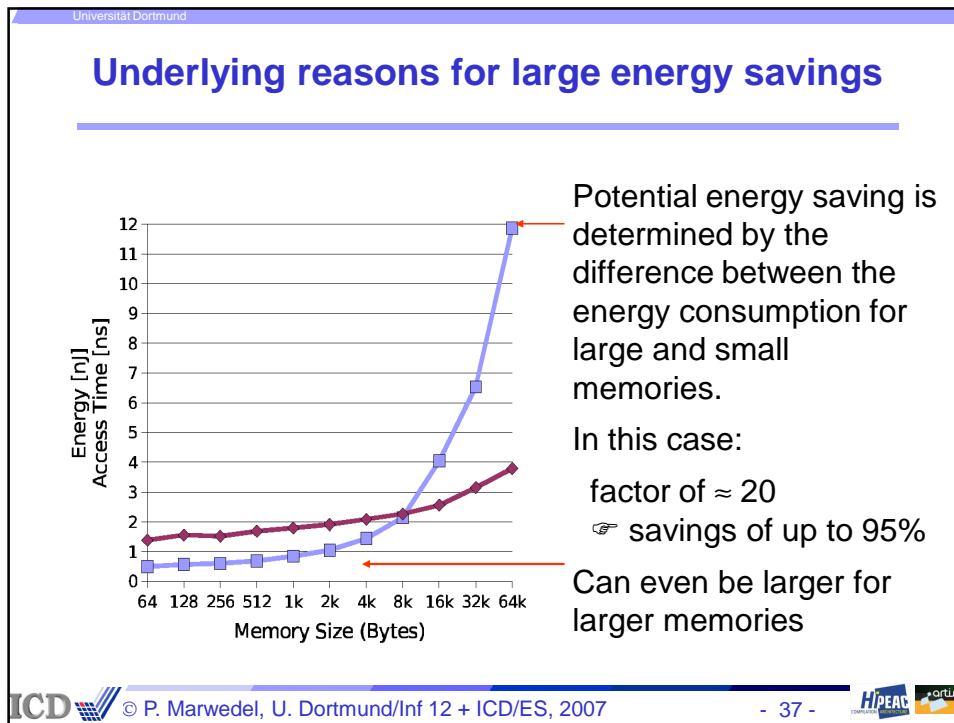
- Consider sets of consecutive BBs as a new kind of basic blocks ("multi blocks")
- Add a constraint preventing the same block from being selected twice:
 
$$x(BB_b) + x(F_i) + \sum_{j \in \text{multiblocks}(b), j \neq x} x(BB_j) \leq 1$$

$$\forall b \in \{\text{blocks}\} \cup \{\text{multi blocks}\}$$

☞ Block  $b$  is either moved individually, as part of a function, as part of one of its enclosing multi-blocks or not at all.

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 34 -








Universitat Dortmund

## Considered partitions

---

# of partitions	number of partitions of size:						
	4k	2k	1k	512	256	128	64
7	0	1	1	1	1	1	2
6	0	1	1	1	1	2	0
5	0	1	1	1	2	0	0
4	0	1	1	2	0	0	0
3	0	1	2	0	0	0	0
2	0	2	0	0	0	0	0
1	1	0	0	0	0	0	0

Example of all considered memory partitions for a total capacity of 4096 bytes

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 39 -  

Universitat Dortmund

## Optimization for multiple memories

---

Minimize  $C = \sum_j (e_j \cdot \sum_i x_{j,i} \cdot n_i)$




With  $e_j$ : energy **required** per access to memory  $j$ ,  
and  $x_{j,i} = 1$  if object  $i$  is mapped to memory  $j$ ,  $=0$  otherwise,  
and  $n_i$ : number of accesses to memory object  $i$ ,  
subject to the constraints:

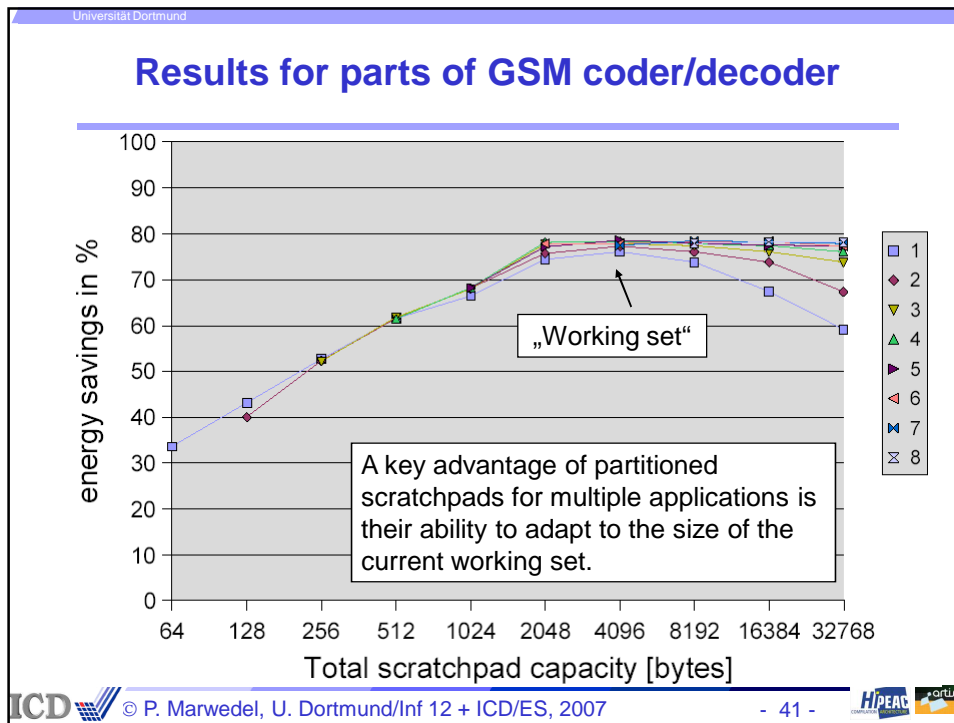
$$\forall j: \sum_i x_{j,i} \cdot S_i \leq SSP_j$$

$$\forall i: \sum_j x_{j,i} = 1$$

With  $S_i$ : size of memory object  $i$ ,  
 $SSP_j$ : size of memory  $j$ .

Main memory  
included as a  
special case of  $j$

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 40 -  



Universität Dortmund

### Summary for lecture 1

- Introduction: Problems with memories:
  - Too slow (☞ memory wall problem)
  - Too energy hungry
  - Not timing predictable
- Scratch pad memories: several advantages
  - Energy-efficient
  - Timing-predictable
  - Small
  - Available (but hardly any tools)
- Non-overlying (“static”) allocation
  - IP models for allocation
    - Functions and variables
    - Fine grained model for basic blocks
    - Partitioning

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 42 - H/PEAC artur



## Compilation for scratch pad memories (2)

### Outline

- Introduction, Motivation
- Scratch-pad memories
  - Characteristics
  - Non-overlaying allocation, partitions, stack, heap, WCET
  - Overlaying (dynamic) allocation
  - Multiple processes & processors
- Caches
  - Predictability, loop transformations
  - Prefetching, layout transformations
- Other levels of the memory hierarchy
  - Array folding
  - ADL-driven compilation for hierarchies
  - Rewriting algorithms for memory hierarchies
- Conclusion, Additional references

Universität Dortmund

## Energy optimization in horizontally partitioned caches

**Energy savings due to two effects**

- Reduction in miss rate
- AccessEnergy (mini cache) < AccessEnergy (main cache)

**Reduction in miss rate**

- Aligned with performance
- Exploited by performance improvement techniques

**Less Energy per Access in mini cache**

- Inverse to performance
- Energy can decrease even if there are more misses*
- Opposite to performance optimization techniques!!*

```

graph TD
    PP[Processor Pipeline] --> MC[Main Cache]
    PP --> MinC[Mini Cache]
    MC --> Mem[Memory]
    MinC --> Mem
    
```

[A. Shrivastava, I. Issenin, N. Dutt: Compilation techniques for energy reduction in horizontally partitioned cache architectures, Intern. Conf. on Compilers, Architectures and Synthesis For Embedded Systems (CASES), 2005, pp. 90-96]

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 45 - H/PEAC

Universität Dortmund

## Simple energy-oriented heuristic are good

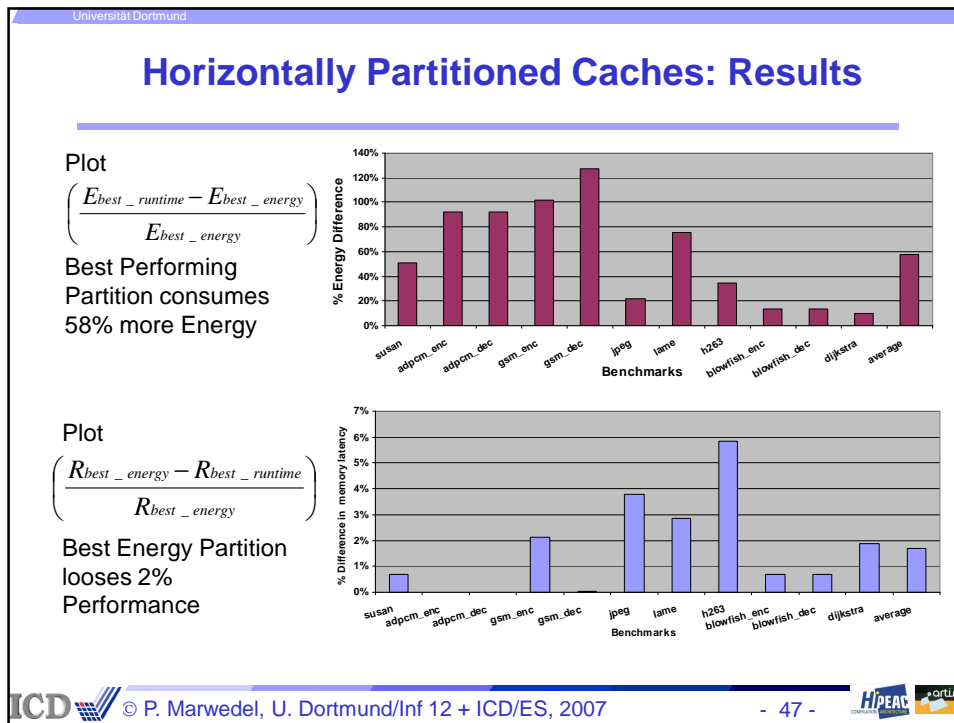
$$\left( \frac{E_{base} - E_{min}}{E_{base}} \right)$$

Benchmark	OPT (%)	OM2N (%)	OMN (%)
susan	~60	~60	~60
adpcm_enc	~45	~45	~45
adpcm_dec	~45	~45	~45
gem_enc	~45	~45	~45
gem_dec	~45	~45	~45
jpeg	~45	~45	~45
lame	~45	~45	~45
h263	~45	~45	~45
blowfish_enc	~45	~45	~45
blowfish_dec	~45	~45	~45
dijkstra	~45	~45	~45
average	~45	~45	~45

- Memory subsystem energy savings achieved by OMN (greedy)
- Base Configuration – All pages are mapped to main cache

OMN achieves 50% memory energy savings

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 46 - H/PEAC



Universität Dortmund

## Stack allocation

1. Run stack size analyzer  
Use **stack as 1 large array** of returned size  
☞ Mapped to SPM if small and frequently accessed  
[approach by Steinke et al. @ Dortmund]
  
2. Potentially **map frequently & less frequently used elements to different memories**  
[approach by Barua et al.]

Considers global & local variables (no instructions)

[O. Avissar, R. Barua, D. Stewart: An Optimal Memory Allocation Scheme for Scratch-Pad-Based Embedded Systems, *Transactions on Embedded Computing Systems*, 2002]

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 48 - H/PEAC artur



Universität Dortmund

## Split stack

---

```

foo()
{
  int a;
  float b;
  ...
}
    
```

Stack in SRAM

Stack in DRAM

☞ 2 stack pointers (overhead)

1. Single stack frame per procedure (fast procedures)
2. Split frames (time consuming procedures)

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 49 -

Universität Dortmund

## IP model for split stack

---

**Cost function (linear):**  
 Run-time = Access time for global variables +  
                     Access time for stack variables

**Constraints (linear):**

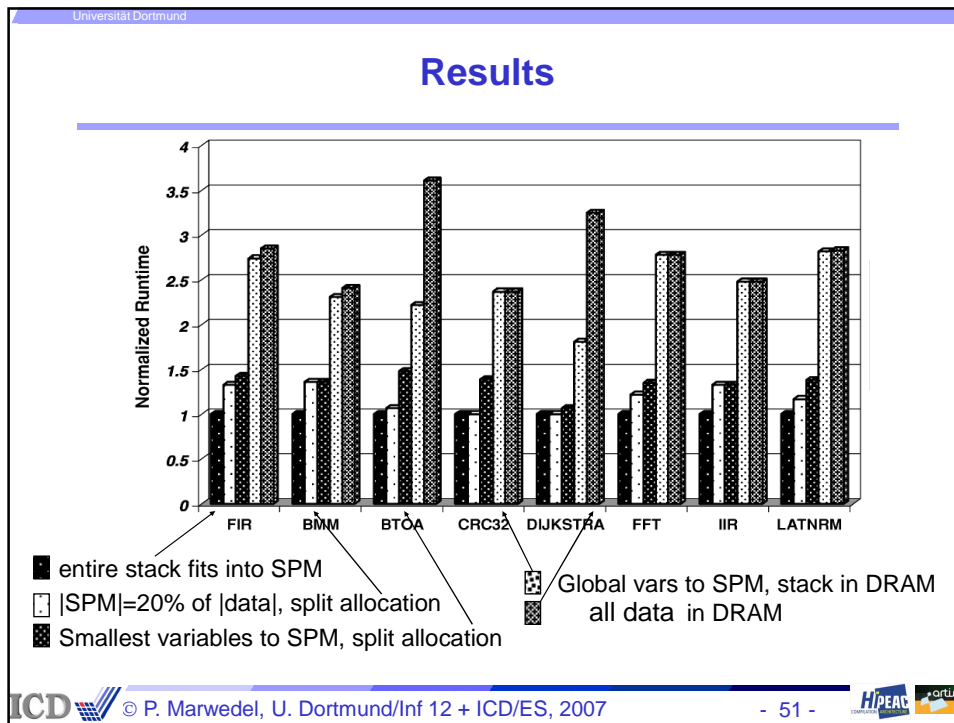
- Each variable has to be mapped to 1 memory
- For all paths in the call graph:  
     The sum of all objects mapped to SPM  
     is  $\leq |SPM|$

```

graph TD
    main((main)) --> foo((foo))
    main --> rom_e((rom e))
    foo --> bar((bar))
    
```

Model actually considers an arbitrary number of memories, each with its own size and access time(s) ☞ "partitioning".

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 50 -



Universität Dortmund

## Heap allocation

1. Run heap size analyzer  
Use **heap as 1 large array** of returned size
2. Potentially **map heap fragments to SPM**  
Problems:
  - Object sizes frequently not known at compile time
  - Avoid illegal references
  - Approaches using a level of indirection suffer from additional overhead

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 52 - H/PEAC artur

Universität Dortmund

## Heap allocation by Barua




---

- Copying between memories is done at *region boundaries* (before & after loops, at procedure entry and exit).
- At each copy point, a *bin* is allocated in SPM for heap data.
- **Objects always have the same address within the *bin*.**
- Modified malloc first checks if space in *bin* is available, otherwise allocates in main RAM
- Modified free deallocates space in *bin*.

Memory offset

	Regions			
	1	2	3	4
256			C	C
512	B	B	B	E
768		A	D	A
1024				

[A. Dominguez, S. Udayakumaran, R. Barua: Heap Data Allocation to Scratch-Pad Memory in Embedded Systems, *Journal of Embedded Computing*, 2005]

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 53 -  



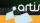
Universität Dortmund

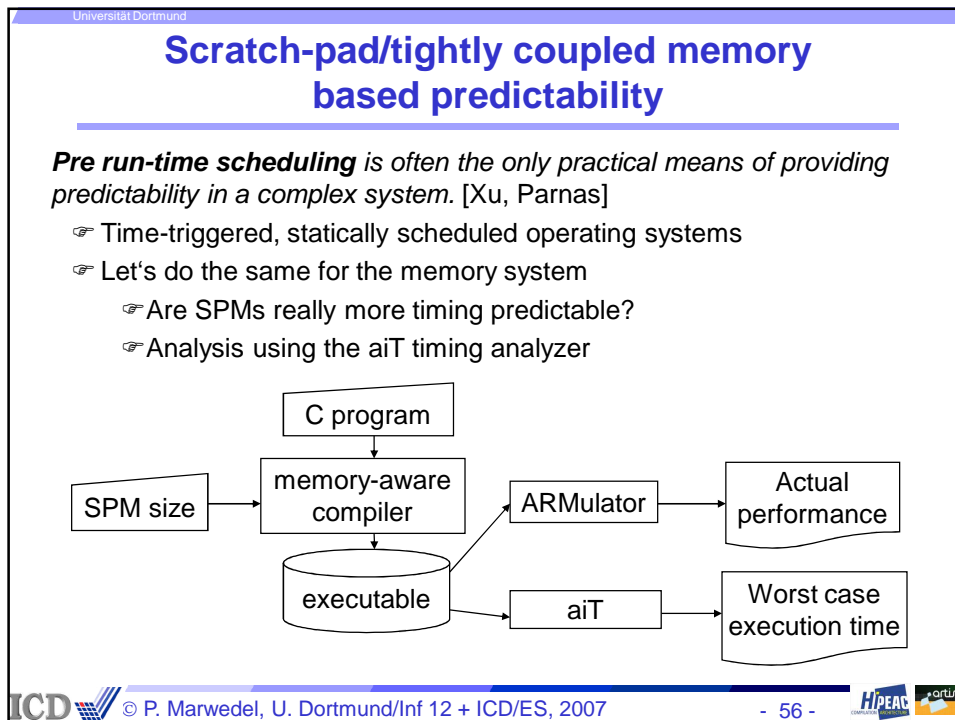
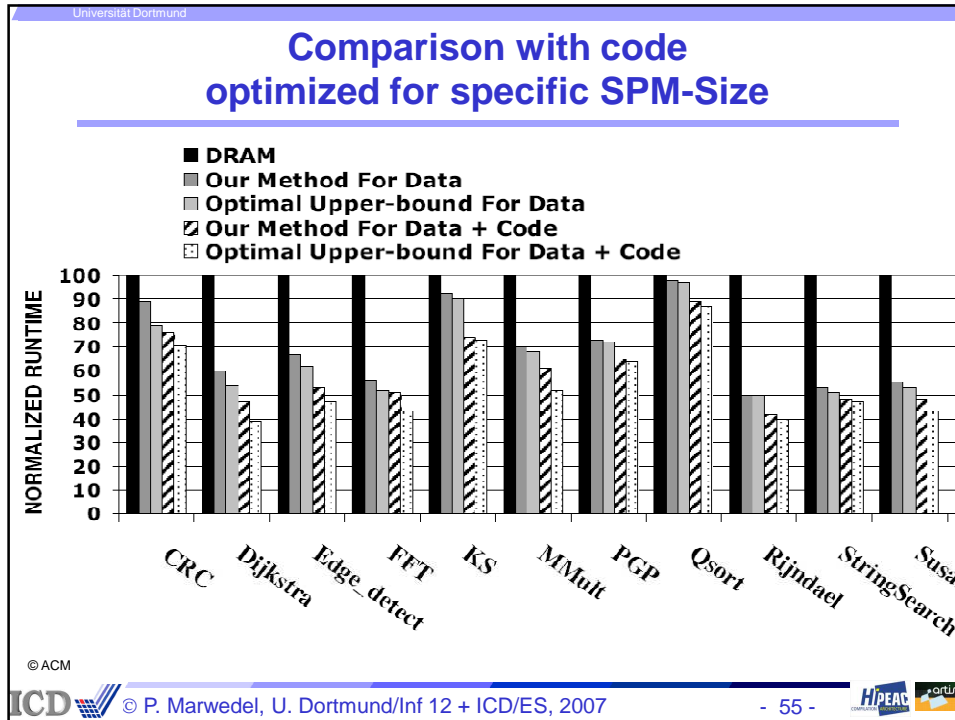
## Considering SPM size at link time

---

- Avoiding executables generated for a specific SPM size
- **Profiler** stores variables sorted by their “frequency per byte” (FPB) in executable
- **Compiler** identifies access to variables with unknown memory allocation by using global symbols.
- After loading the program, a custom **installer**
  - Reads SPM size of current architecture
  - Decides (using FPBs) which variables to put into SPM
  - Patches addresses of variables (incl. stack variables)
  - Tries to use SPM space not used on a calling path
  - Also moves code into SPM

[N. Nguyen, A. Dominguez, R. Barua: Memory Allocation for Embedded Systems with a Compile-Time-Unknown Scratch-Pad Size, *Intern. conf. on Compilers, architectures and synthesis for embedded systems (CASES)*, 2005, p. 115-125]

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 54 -  






Universitat Dortmund

## Architectures considered

---

ARM7TDMI with 3 different memory architectures:

- 1. Main memory**  
 LDR-cycles: (CPU,IF,DF)=(3,2,2)  
 STR-cycles: (2,2,2)  
 \* = (1,2,0)
- 2. Main memory + unified cache**  
 LDR-cycles: (CPU,IF,DF)=(3,12,6)  
 STR-cycles: (2,12,3)  
 \* = (1,12,0)
- 3. Main memory + scratch pad**  
 LDR-cycles: (CPU,IF,DF)=(3,0,2)  
 STR-cycles: (2,0,0)  
 \* = (1,0,0)

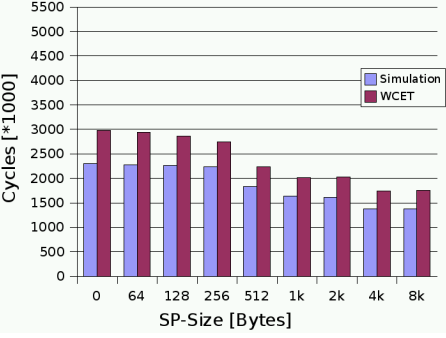
ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 57 -  

Universitat Dortmund

## Results for G.721

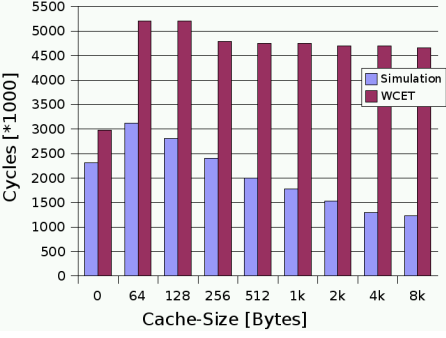
---

**Using Scratchpad:**



SP-Size [Bytes]	Simulation Cycles [*1000]	WCET Cycles [*1000]
0	2200	3000
64	2200	2900
128	2200	2800
256	2200	2700
512	2200	2400
1k	2200	2100
2k	2200	2000
4k	2200	1800
8k	2200	1700




**Using Unified Cache:**



Cache-Size [Bytes]	Simulation Cycles [*1000]	WCET Cycles [*1000]
0	2200	3000
64	3000	5000
128	2800	5000
256	2400	4800
512	2000	4700
1k	1800	4700
2k	1500	4700
4k	1300	4700
8k	1200	4700

**Yes, they are clearly more timing predictable!**




- L. Wehmeyer, P. Marwedel: Influence of Onchip Scratchpad Memories on WCET: *4th Intl Workshop on worst-case execution time analysis, (WCET)*, 2004
- L. Wehmeyer, P. Marwedel: Influence of Memory Hierarchies on Predictability for Time Constrained Embedded Software, *Design Automation and Test in Europe (DATE)*, 2005

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 58 -  

Universität Dortmund

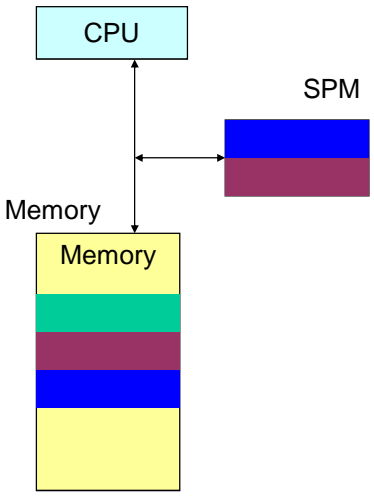
## Outline

- Introduction, Motivation
- Scratch-pad memories
  - Characteristics
  - Non-overlying allocation, partitions, stack, heap, WCET
  - ➡ - Overlaying (dynamic) allocation
  - Multiple processes & processors
- Caches
  - Predictability, loop transformations
  - Prefetching, layout transformations
- Other levels of the memory hierarchy
  - Array folding
  - ADL-driven compilation for hierarchies
  - Rewriting algorithms for memory hierarchies
- Conclusion, Additional references




ICD  © P. Marwedel, U. Dortmund/Inf. 12 + ICD/ES, 2007 - 59  

Universität Dortmund

## Non-overlying allocation problematic for multiple hot spots ➡ Overlaying allocation



- Effectively results in a kind of **compiler-controlled overlays** for SPM
- Address assignment within SPM required

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 60 -  

Universität Dortmund

## Overlaying allocation by Barua et al. (1)

Potential copying @ control points (major locality changes):

- Beginning and end of procedures
- Before and after loops
- Before and after if statements
- Beginning and end of then and else parts
- Before and after switch statements and their cases

Number control points!

```

main(){
  if(...) {proc-D()} else {...}
  proc-A()
  proc-B()
}
proc-A() {
  proc-C()
}
proc-B {
  proc-C()
  while(...) {Y=...}
}
proc-C() {X=...}
proc-D() {...}
  
```

[S. Udayakumaran, A. Dominguez, R. Barua: Dynamic Allocation for Scratch-Pad Memory using Compile-Time Decisions, *ACM Transactions in Embedded Computing Systems*, Vol. V, 2006, Pages 472 - 511]

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 61 -

Universität Dortmund

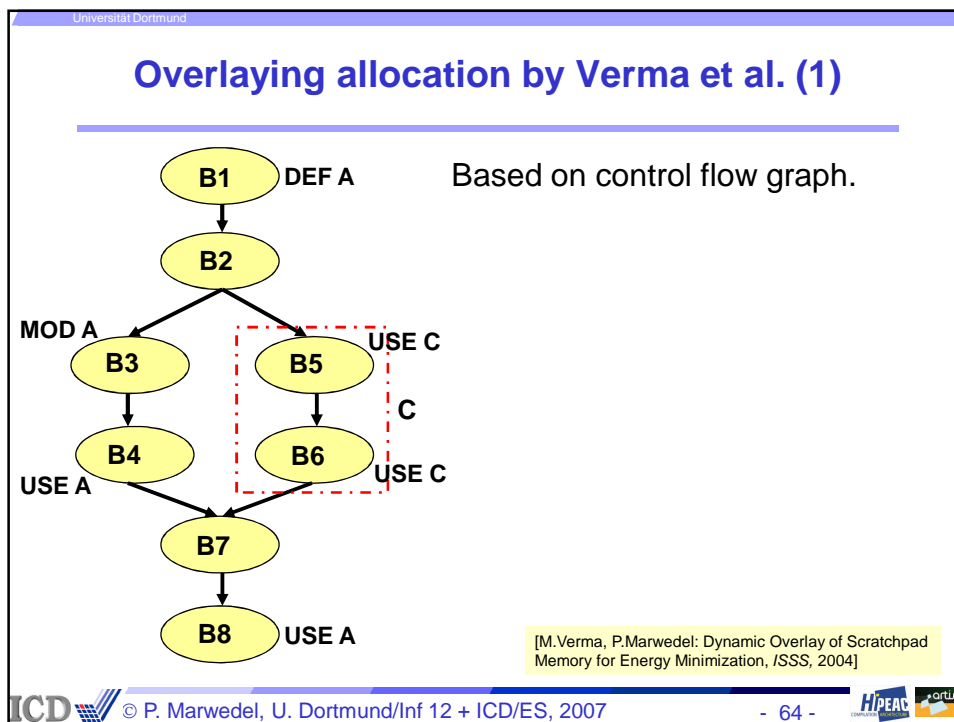
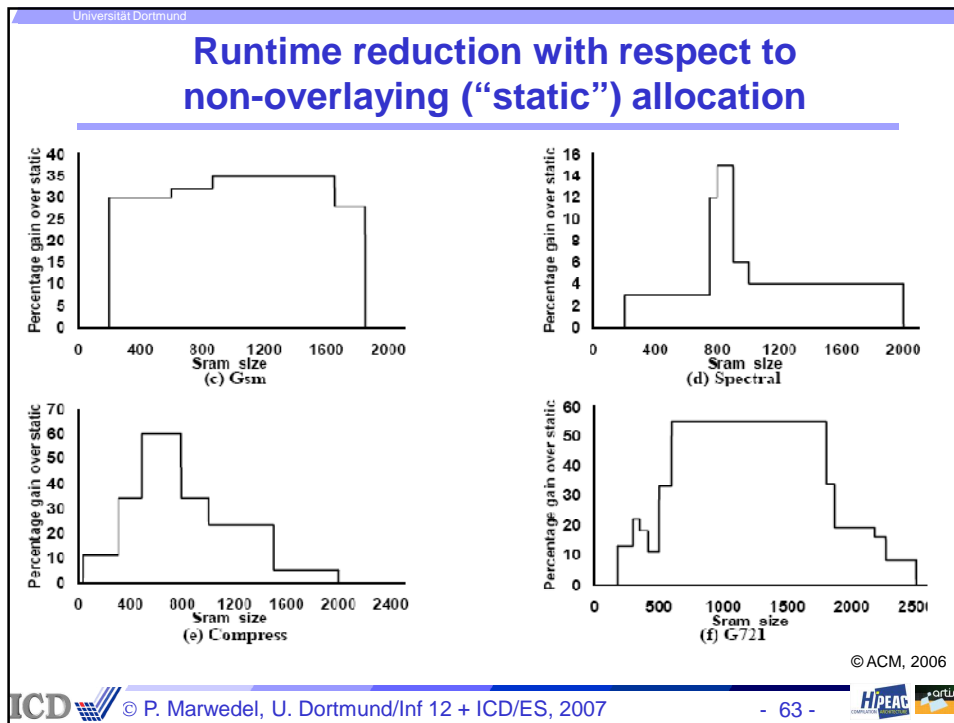
## Overlaying allocation by Barua et al. (2)

- At 1st control point: Initialize SPM with variables accessed in first region in decreasing order of freq-per-byte
- Traverse control points according to control point #

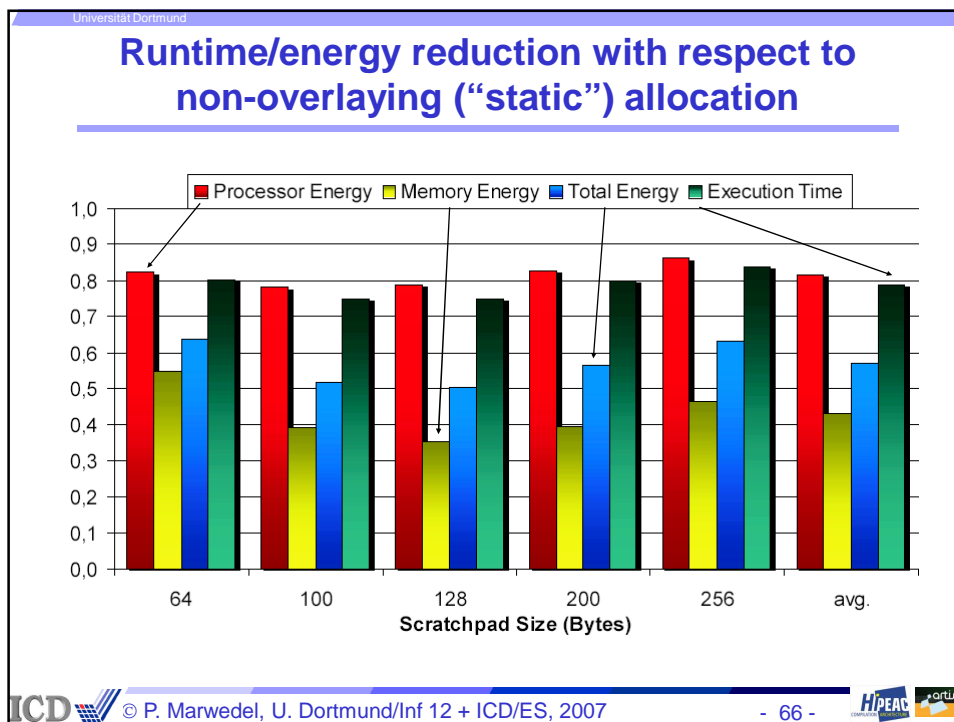
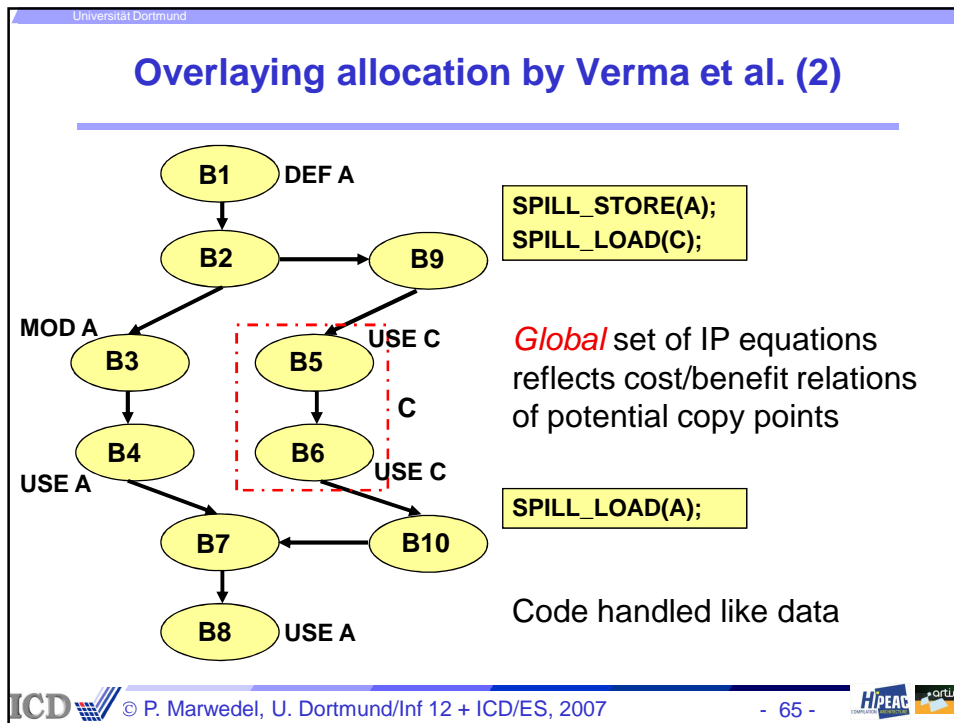
At each control point:

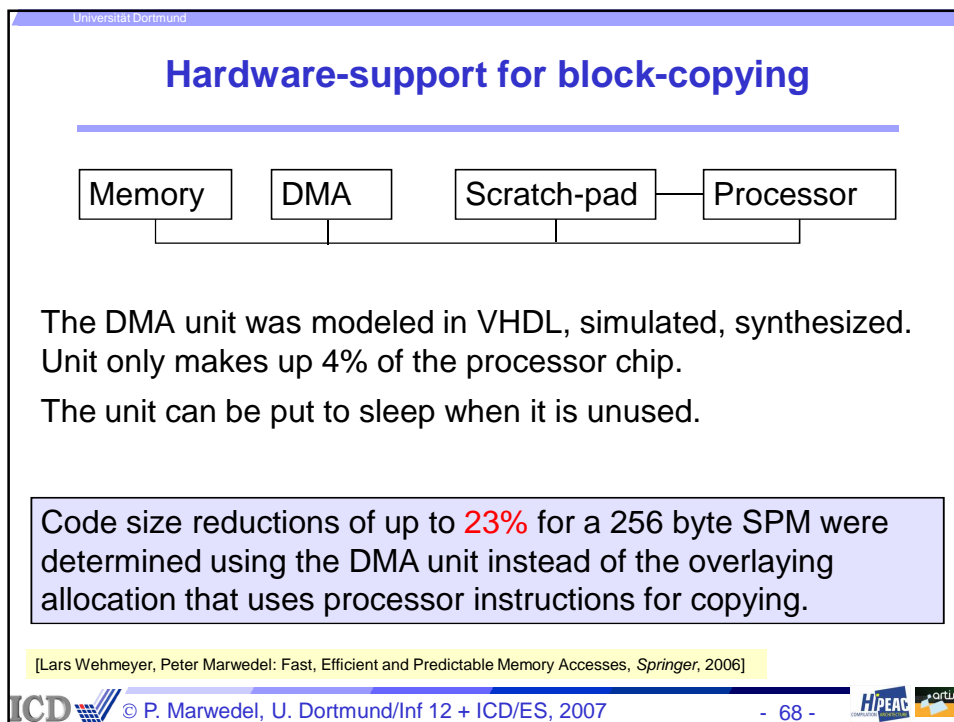
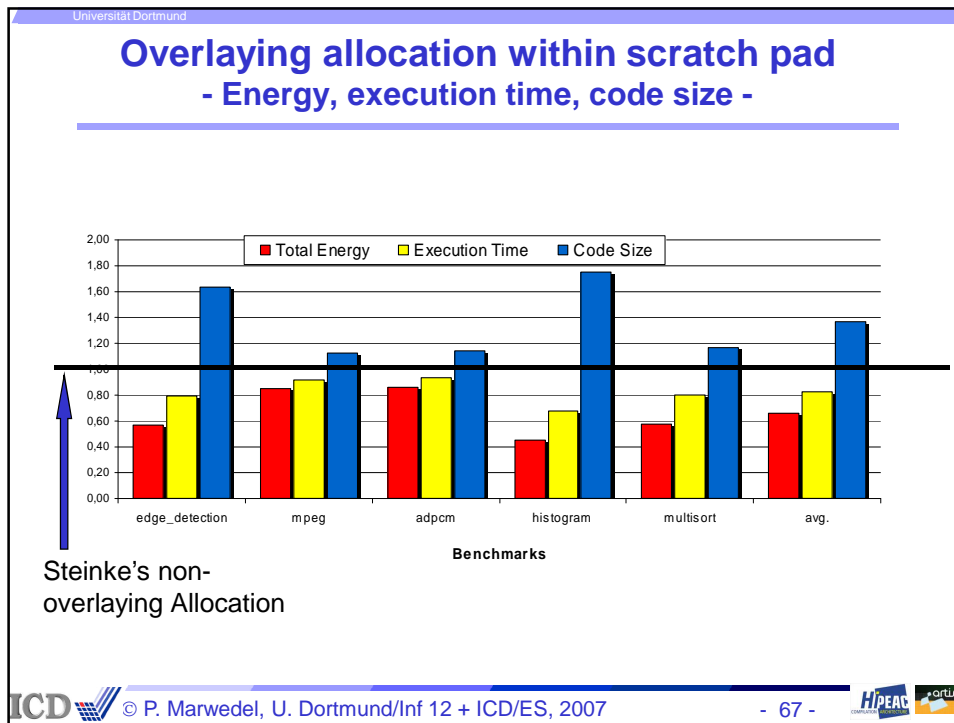
- Compute variables that potentially should be swapped in, based on # of accesses per byte
- Compute variables that should be swapped out
- Compute variables that should remain
- Compute addresses of variables to be moved in
- Generate swap code

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 62 -









Universitat Dortmund

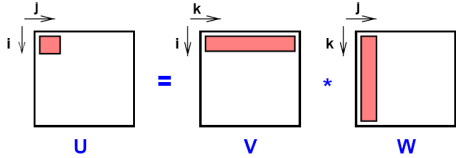
## References to large arrays (1)

### - Regular accesses -

```

for (i=0; i<n; i++)
  for (j=0; j<n; j++)
    for (k=0; k<n; k++)
      U[i][j]=U[i][j] + V[i][k] * W[k][j]

```



**Tiling**

```

for (it=0; it<n; it=it+Sb)
  {read_tile V[it:it+Sb-1, 1:n]
  for (jt=0; jt<n; jt=jt+Sb)
    {read_tile U[it:it+Sb-1, jt:jt+Sb-1]
    read_tile W[1:n, jt:jt+Sb-1]
    U[it:it+Sb-1, jt:jt+Sb-1]=U[it:it+Sb-1, jt:jt+Sb-1]
      + V[it:it+Sb-1, 1:n]
      * W[1:n, jt:jt+Sb-1]
    write_tile U[it:it+Sb-1, jt:jt+Sb-1]
  }
}

```

[M. Kandemir, J. Ramanujam, M. J. Irwin, N. Vijaykrishnan, I. Kadayif, A. Parikh: Dynamic Management of Scratch-Pad Memory Space, *DAC*, 2001, pp. 690-695]

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 69 -

Universitat Dortmund

## References to large arrays

### - Irregular accesses -

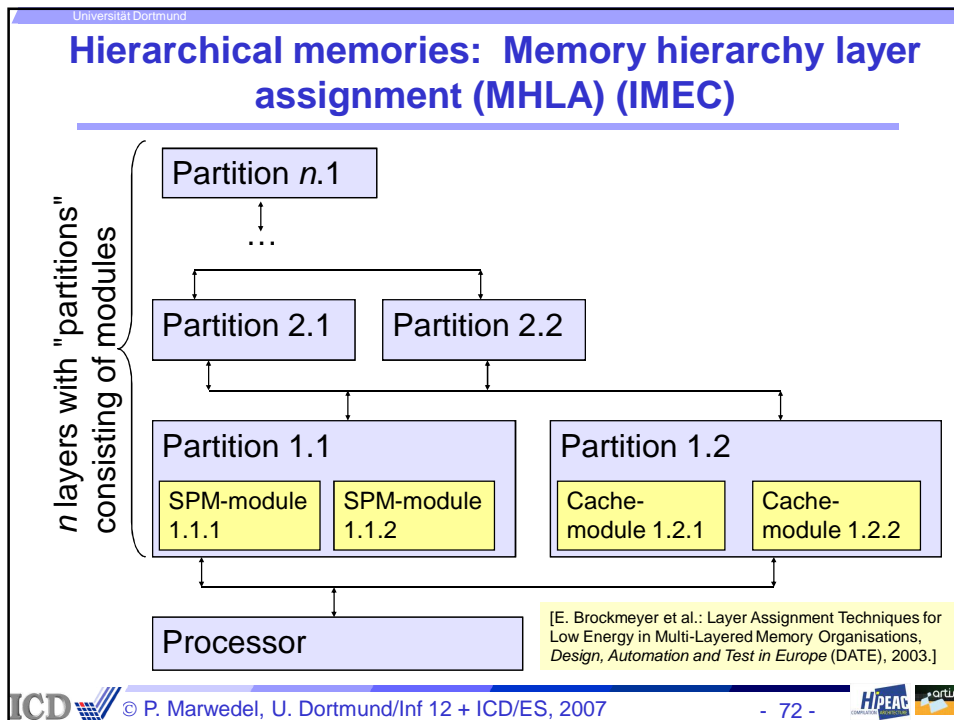
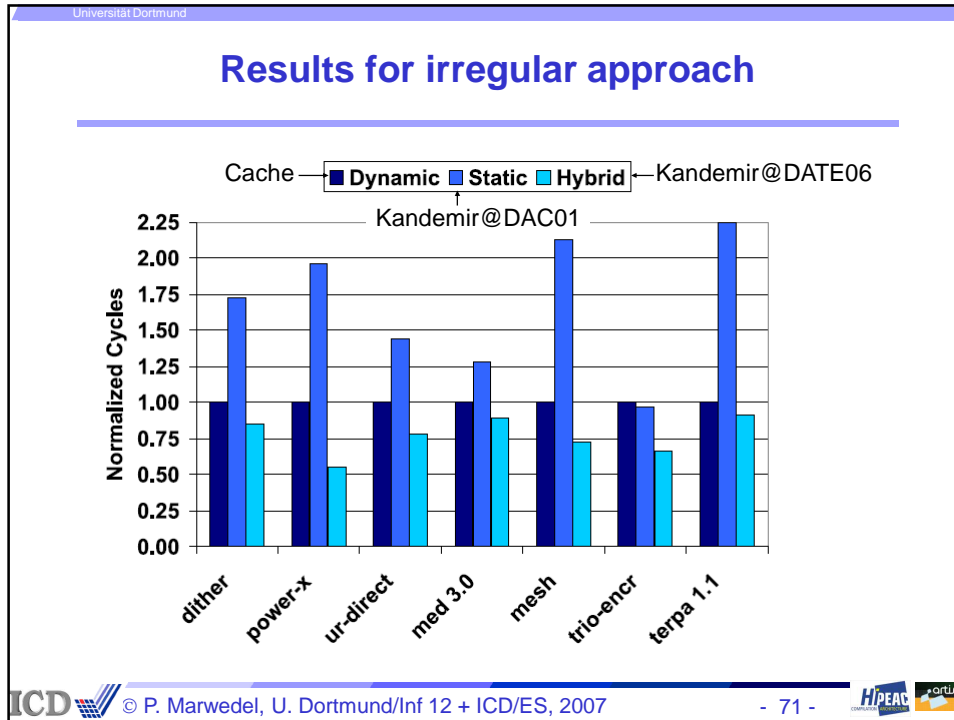
```

for each loop nest L in program P {
  apply loop tiling to L based on the access patterns of
  regular array references;
  for each assignment to index array X
    update the block minimum and maximum values of X;
  compute the set of array elements that are irregularly
  referenced in the current inter-tile iteration;
  compare the memory access costs for using
  and not using SPM;
  if (using SPM is beneficial)
    execute the intra-tile loop iterations by using the SPM
  else
    execute the intra-tile loop iterations by not
    using the SPM
}

```

[G. Chen, O. Ozturk, M. Kandemir, M. Karakoy: Dynamic Scratch-Pad Memory Management for Irregular Array Access Patterns, *DATE*, 2006]

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 70 -



Universität Dortmund

## Memory hierarchy layer assignment (MHLA) - Copy candidates -

```
int A[250]
for (i=0; i<10; i++)
  for (j=0; j<10; j++)
    for (k=0; k<10; k++)
      for (l=0; l<10; l++)
        f(A[j*10+l])
size=0; reads(A)=10000
```

```
int A[250]
for (i=0; i<10; i++)
  for (j=0; j<10; j++)
    {A'[0..9]=A[j*10..j*10+9];
     for (k=0; k<10; k++)
       for (l=0; l<10; l++)
         f(A'[l])}
size=10; reads(A)=1000
```

Copy candidate  
A', A" in small memory

```
int A[250]
for (i=0; i<10; i++)
  {A'[0..99]=A[0..99];
   for (j=0; j<10; j++)
     for (k=0; k<10; k++)
       for (l=0; l<10; l++)
         f(A'[j*10+l])}
size=100; reads(A)=1000
```

```
int A[250]
A'[0..99]=A[0..99];
for (i=0; i<10; i++)
  for (j=0; j<10; j++)
    for (k=0; k<10; k++)
      for (l=0; l<10; l++)
        f(A'[j*10+l])
size=100; reads(A)=100
```

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 73 -

Universität Dortmund

## Memory hierarchy layer assignment (MHLA) - Goal -

**Goal:** For each variable: find permanent layer, partition and module & select copy candidates such that energy is minimized.

### Conflicts between variables

Nr misses = 3000

Nr misses = 2750

Nr misses = 5250

Nr misses = 5250

[E. Brockmeyer et al.: Layer Assignment Techniques for Low Energy in Multi-Layered Memory Organisations, *Design, Automation and Test in Europe (DATE)*, 2003.]

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 74 -

Universität Dortmund

## Memory hierarchy layer assignment (MHLA) - Approach -

**Approach:**

- start with initial variable allocation
- incrementally improve initial solution such that total energy is minimized.

Current assignment

Next assignment

More general hardware architecture than the Dortmund approach, but no global optimization.

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 75 - H/PEAC artur

Universität Dortmund

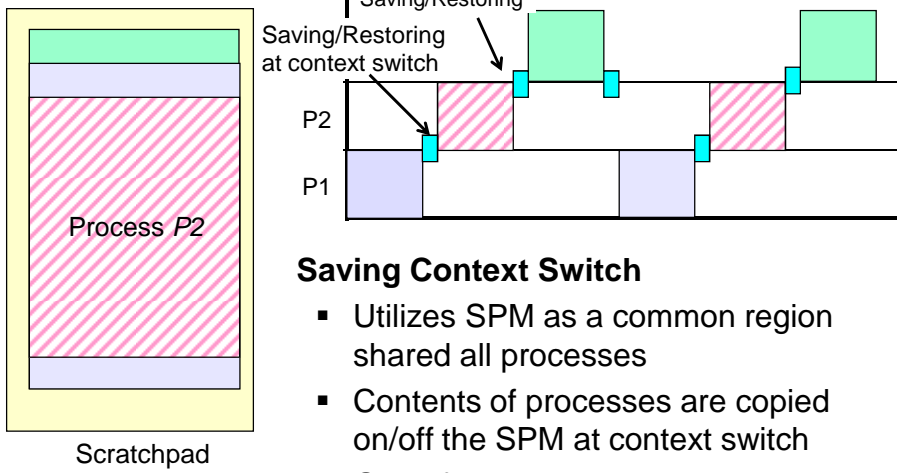
## Outline

- Introduction, Motivation
- Scratch-pad memories
  - Characteristics
  - Non-overlying allocation, partitions, stack, heap, WCET
  - Overlaying (dynamic) allocation
- ➡ ▪ Multiple processes & processors
- Caches
  - Predictability, loop transformations
  - Prefetching, layout transformations
- Other levels of the memory hierarchy
  - Array folding
  - ADL-driven compilation for hierarchies
  - Rewriting algorithms for memory hierarchies
- Conclusion, Additional references

ICD © P. Marwedel, U. Dortmund/Inf. 12 + ICD/ES, 2007 - 76 - H/PEAC artur

Universität Dortmund

## Static set of multiple applications: Saving/Restoring Context Switch



The diagram illustrates the context switch mechanism for a static set of multiple applications. On the left, a vertical stack of three colored rectangles (green, light blue, and pink) represents the memory layout of Process P2, with the pink region labeled 'Scratchpad'. On the right, a horizontal timeline shows the execution of three processes: P1 (blue), P2 (pink), and P3 (green). The timeline is divided into segments by vertical lines. Small cyan squares at the boundaries of these segments indicate the points where context is saved or restored. Arrows labeled 'Saving/Restoring at context switch' point to these cyan squares. The processes are shown in a non-overlapping, sequential manner, indicating that only one process is active at any given time.

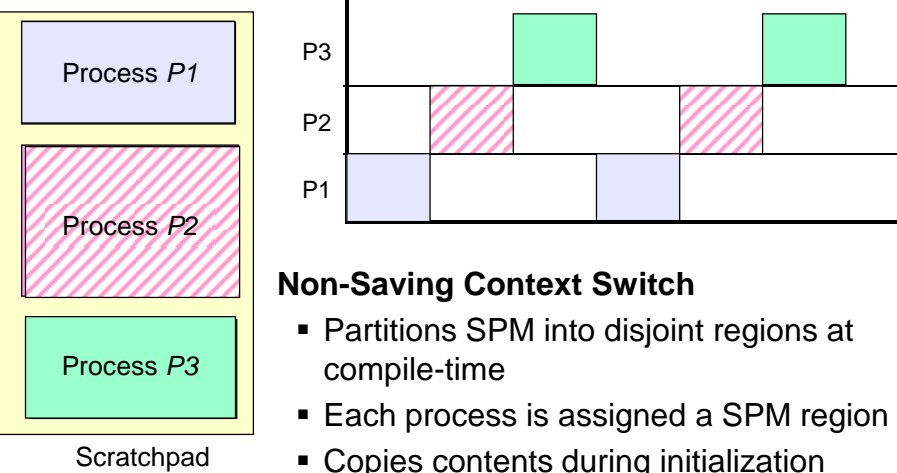
**Saving Context Switch**

- Utilizes SPM as a common region shared all processes
- Contents of processes are copied on/off the SPM at context switch
- Good for small scratchpads

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 77 - H/PEAC

Universität Dortmund

## Non-Saving Context Switch

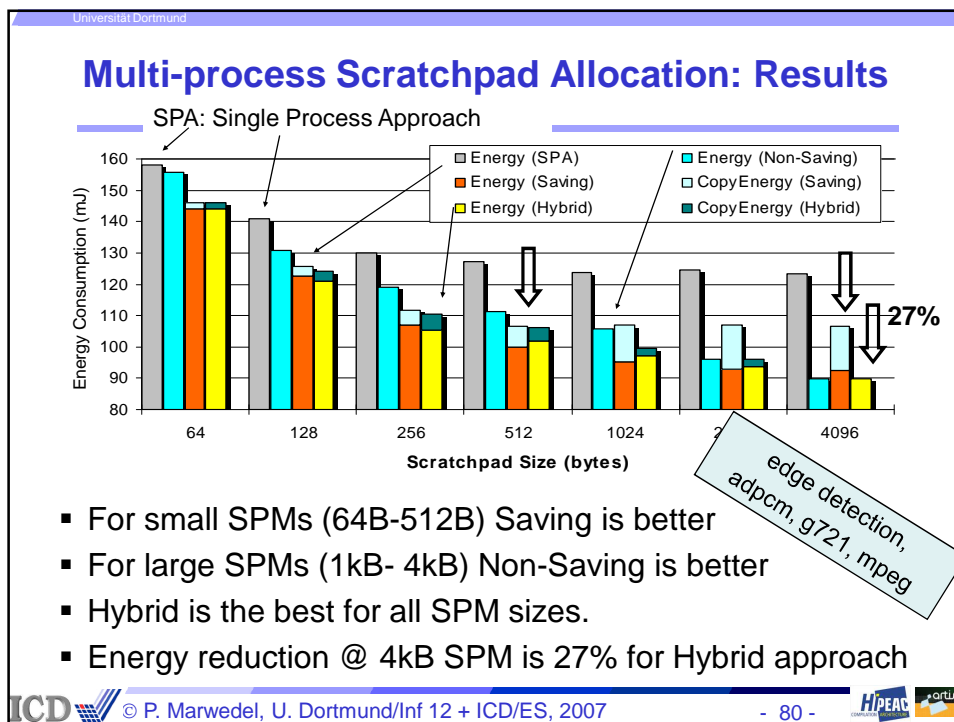
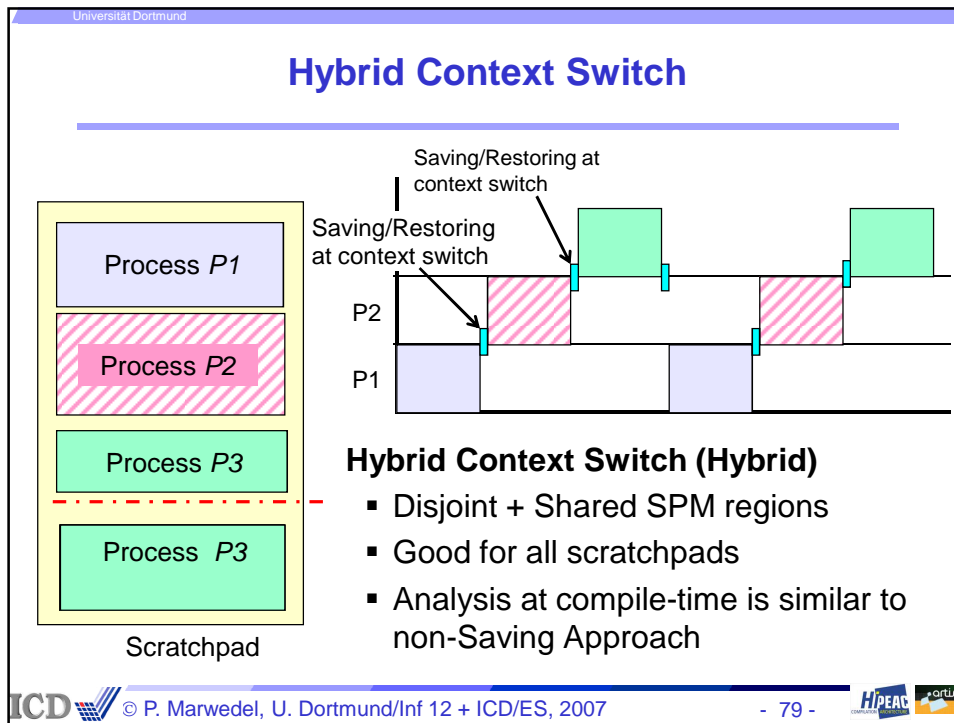


The diagram illustrates the context switch mechanism for a non-saving context switch. On the left, a vertical stack of three colored rectangles (blue, pink, and green) represents the memory layout of three processes: Process P1 (blue), Process P2 (pink), and Process P3 (green), with the pink region labeled 'Scratchpad'. On the right, a horizontal timeline shows the execution of three processes: P1 (blue), P2 (pink), and P3 (green). The timeline is divided into segments by vertical lines. The processes are shown in a non-overlapping, sequential manner, indicating that only one process is active at any given time. Unlike the saving/restoring mechanism, there are no cyan squares at the boundaries of the segments, indicating that context is not saved or restored during the switch.

**Non-Saving Context Switch**

- Partitions SPM into disjoint regions at compile-time
- Each process is assigned a SPM region
- Copies contents during initialization
- Good for large scratchpads

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 78 - H/PEAC





## Summary lecture 2

- Non-Overlaying allocation
  - Stack allocation
  - Heap allocation
  - Timing predictability
- Overlaying SPM allocation
  - Single process
    - Barua's call graph approach
    - Verma's CDFG-based approach
    - Kandemir's tiling approach
    - IMEC's multiple levels
  - Multiple processes
    - Dortmund approach: saving, non-saving, hybrid



**Compilation for**  
**- Scratch pad memories (3)**  
**- Caches**

Universität Dortmund

## Outline

---

- Introduction, Motivation
- Scratch-pad memories
  - Characteristics
  - Non-overlying allocation, partitions, stack, heap, WCET
  - Overlaying (dynamic) allocation
  - Multiple processes & processors
- Caches
  - Predictability, loop transformations
  - Prefetching, layout transformations
- Other levels of the memory hierarchy
  - Array folding
  - ADL-driven compilation for hierarchies
  - Rewriting algorithms for memory hierarchies
- Conclusion, Additional references

ICD © P. Marwedel, U. Dortmund/Inf. 12 + ICD/ES, 2007 - 83 -

Universität Dortmund

## Dynamic set of multiple applications

---

Compile-time partitioning of SPM no longer feasible

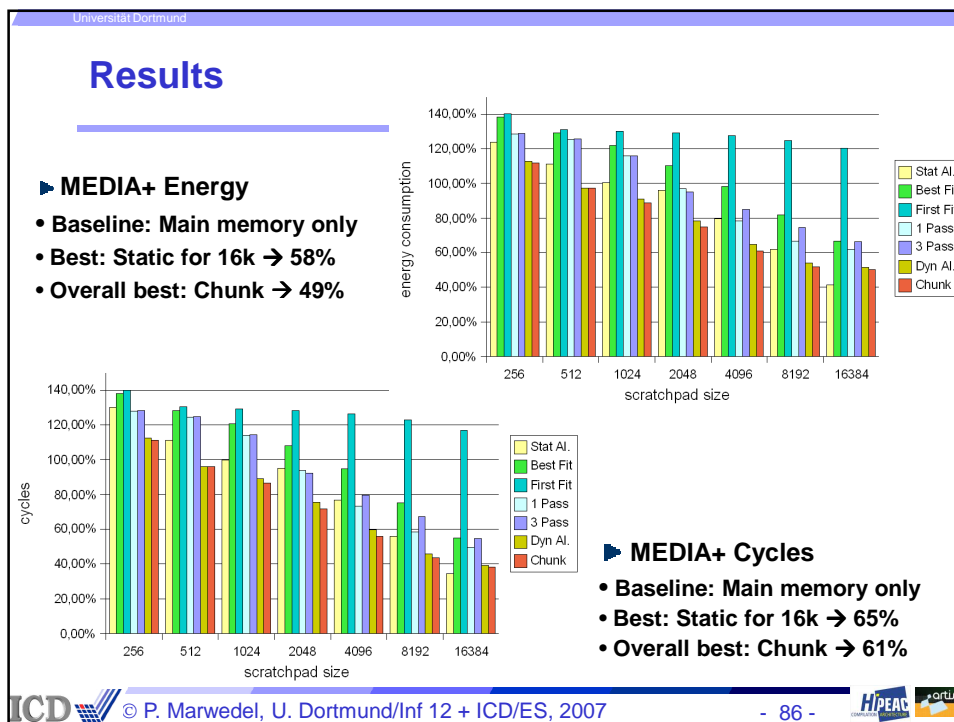
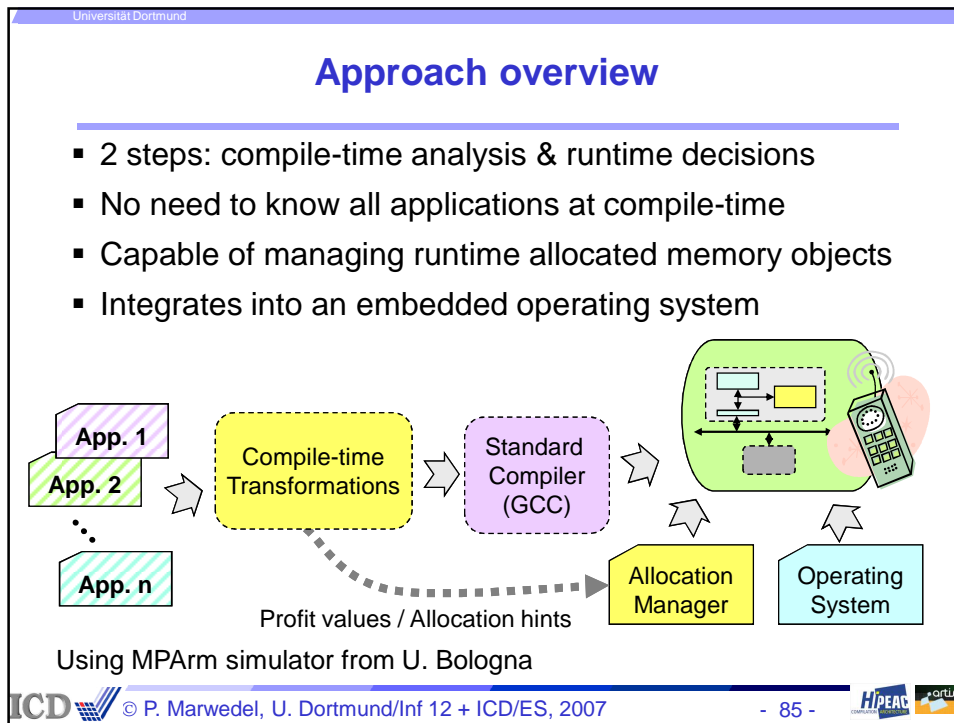
☞ Introduction of SPM-manager

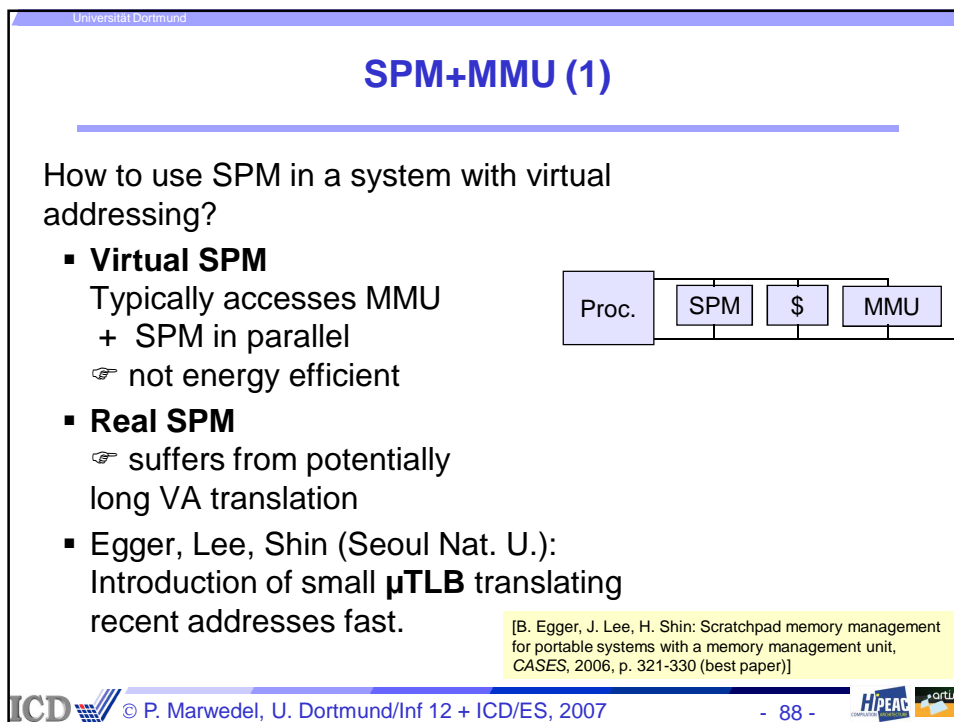
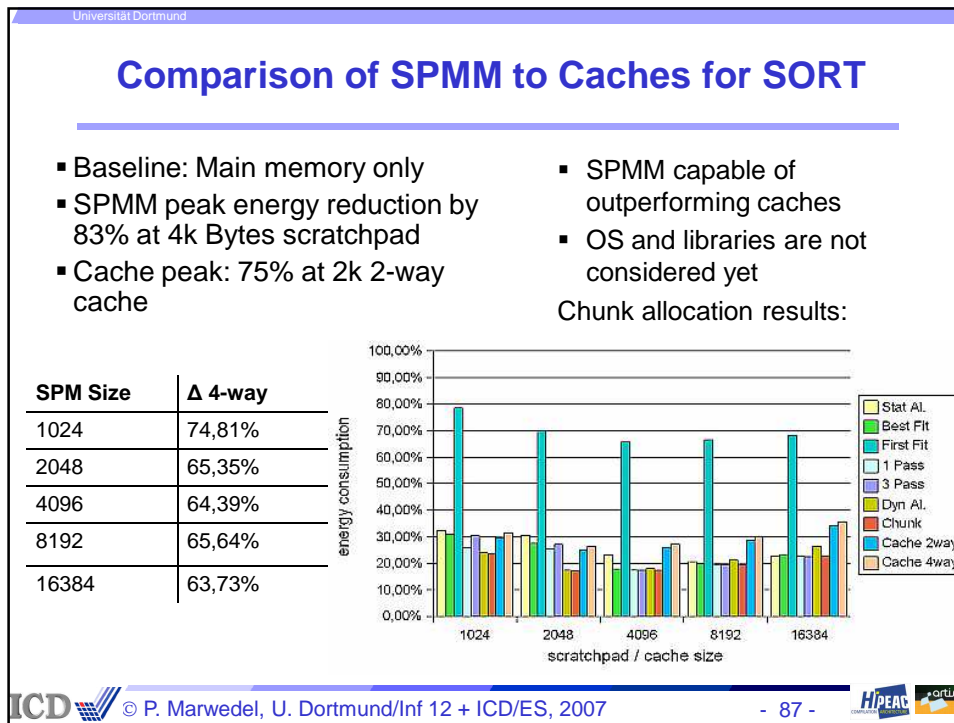
- Runtime decisions, but compile-time supported

Address space:

[R. Pyka, Ch. Faßbach, M. Verma, H. Falk, P. Marwedel: Operating system integrated energy aware scratchpad allocation strategies for multi-process applications, *SCOPES*, 2007]

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 84 -





### SPM+MMU (2)

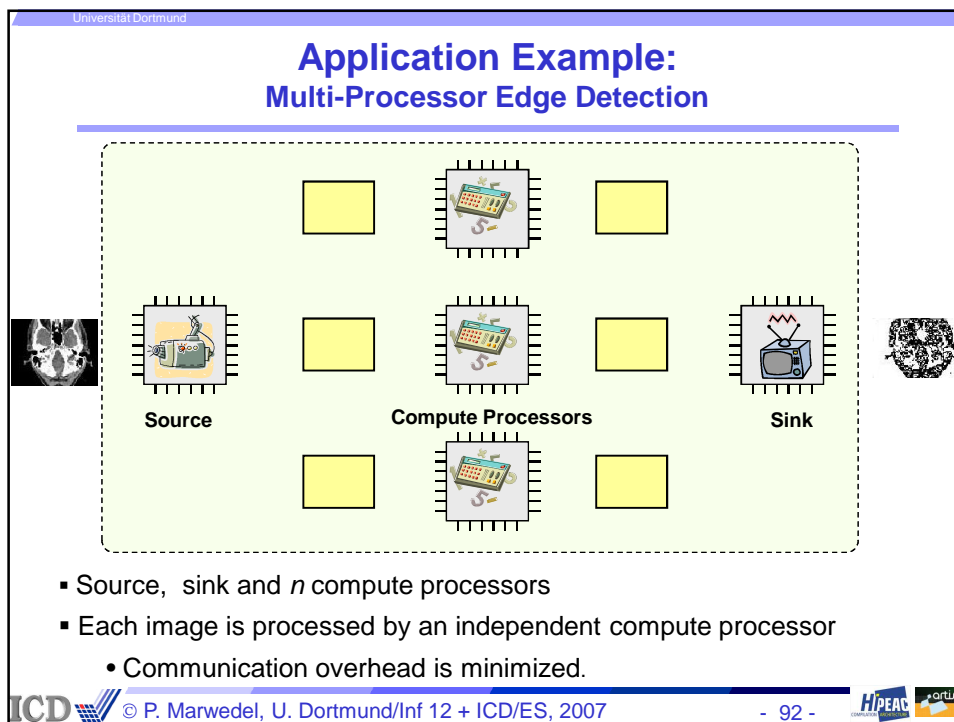
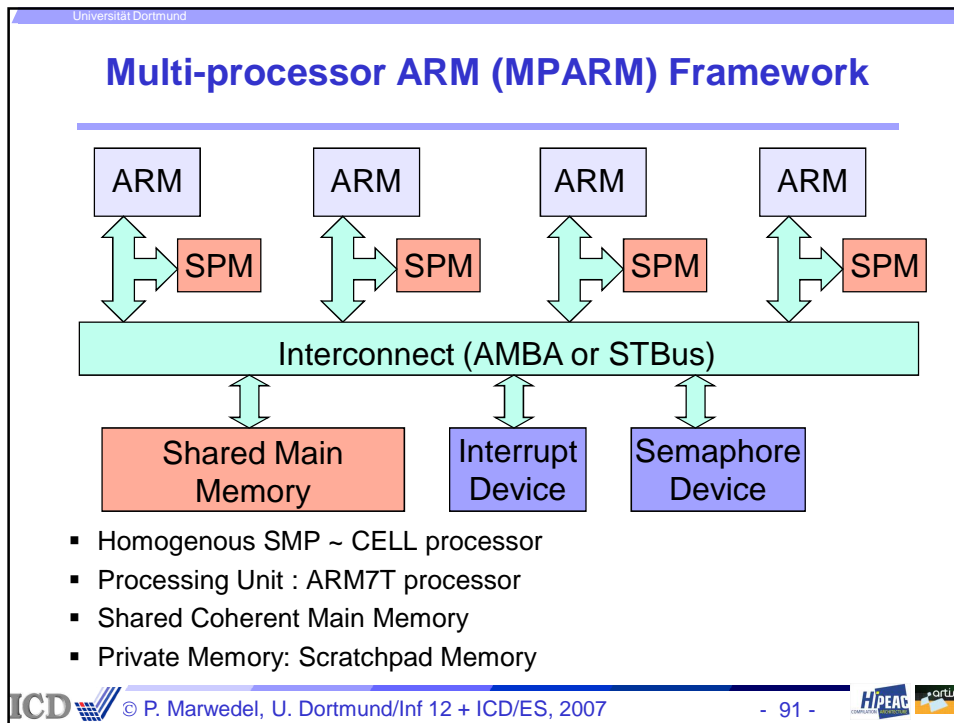
- $\mu$ TLB generates physical address in 1 cycle
- if address corresponds to SPM, it is used
- otherwise, mini-cache is accessed
- Mini-cache provides reasonable performance for non-optimized code
- $\mu$ TLB miss triggers main TLB/MMU
- SPM is used only for instructions
- instructions are stored in pages
- pages are classified as cacheable, non-cacheable, and "pageable" (= suitable for SPM)

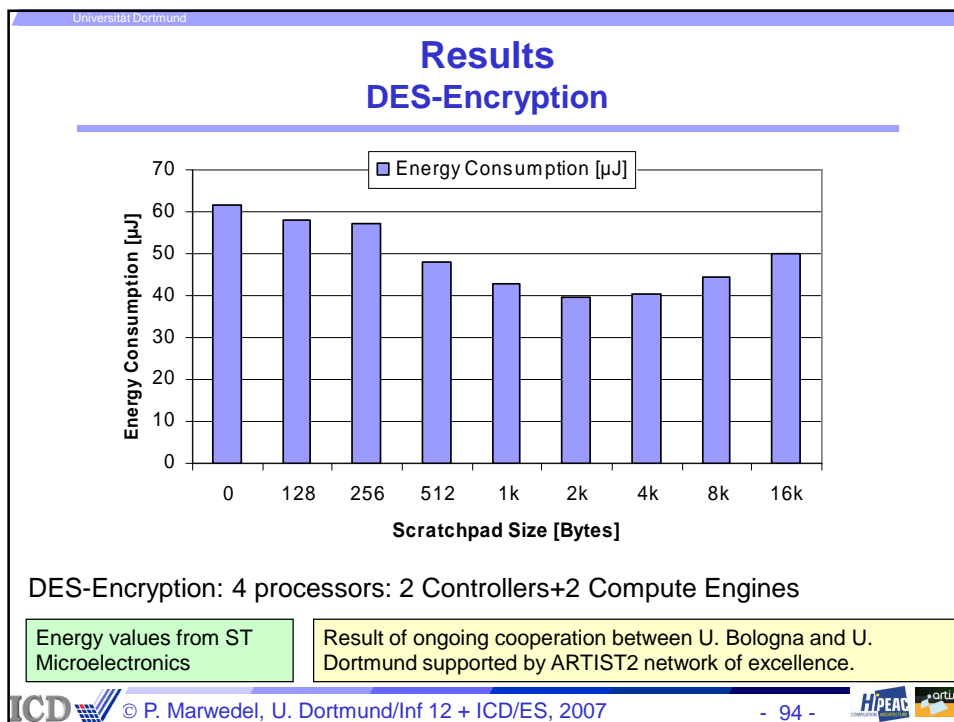
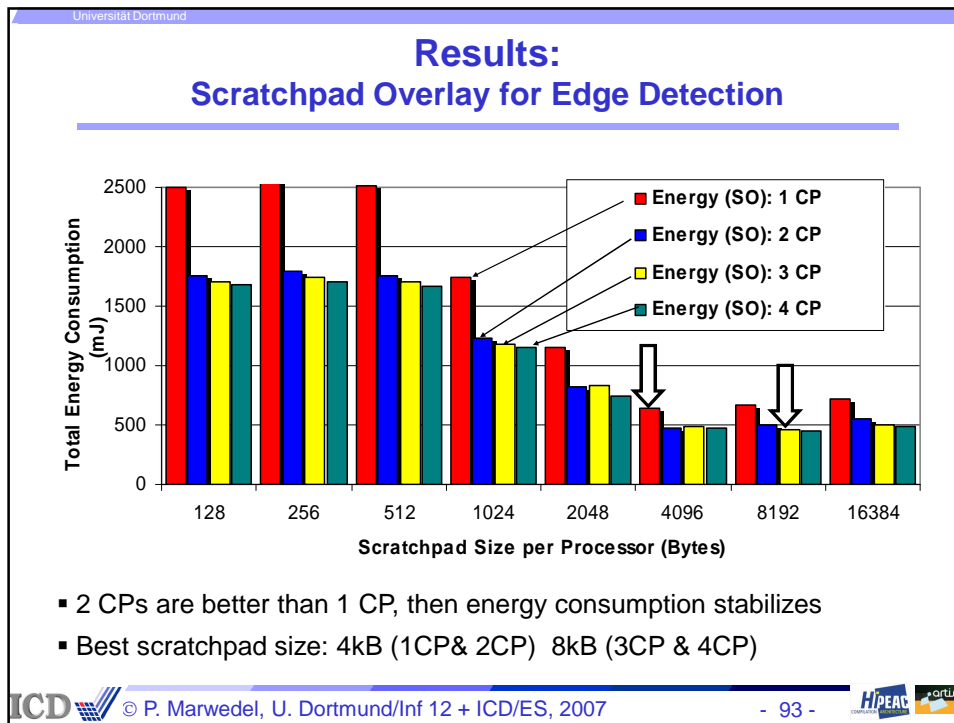
ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 89 -

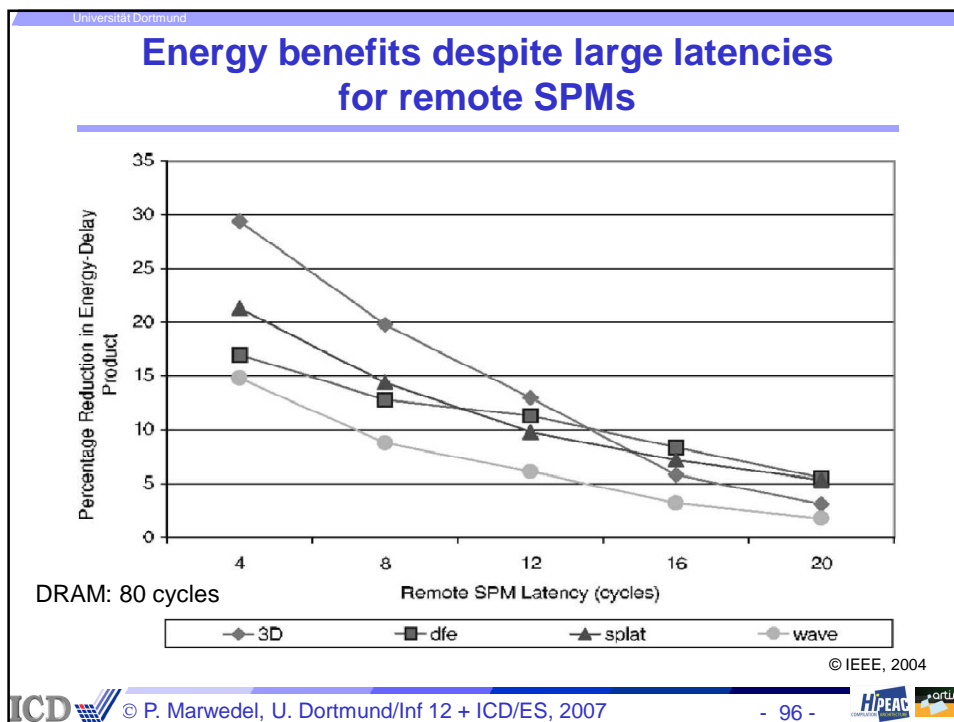
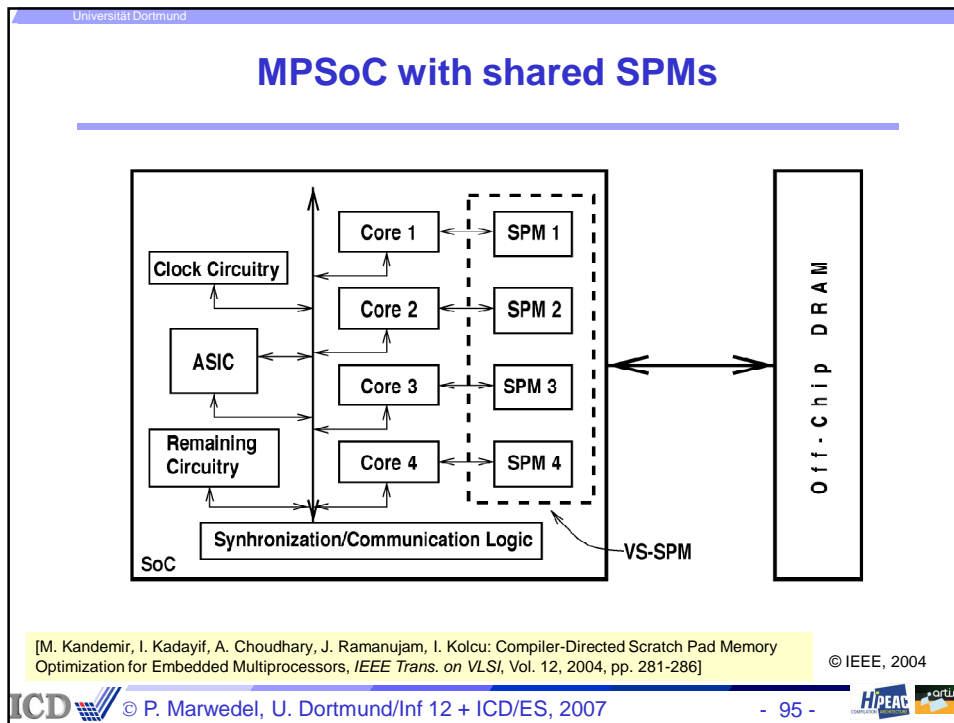
### SPM+MMU (3)

- Application binaries are modified: frequently executed code put into pageable pages.
- Initially, page-table entries for pageable code are marked invalid
- If invalid page is accessed, a *page table exception* invokes SPM manager (SPMM).
- SPMM allocates space in SPM and sets page table entry
- If SPMM detects more requests than fit into SPM, SPM eviction is started
- Compiler does not need to know SPM size

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 90 -












Universität Dortmund

## Extensions

- Using DRAM
- Applications to Flash memory (copy code or execute in place):

} PhD thesis of Lars Wehmeyer  
according to own experiments: very much parameter dependent




- Trying to imitate advantages of SPM with caches: partitioned caches, etc.

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 97 -  

Universität Dortmund

## Outline

- Introduction, Motivation
- Scratch-pad memories
  - Characteristics
  - Non-overlying allocation, partitions, stack, heap, WCET
  - Overlaying (dynamic) allocation
  - Multiple processes & processors
- ➔ Caches
  - Predictability, loop transformations
  - Prefetching, layout transformations
- Other levels of the memory hierarchy
  - Array folding
  - ADL-driven compilation for hierarchies
  - Rewriting algorithms for memory hierarchies
- Conclusion, Additional references

ICD  © P. Marwedel, U. Dortmund/Inf. 12 + ICD/ES, 2007 - 98 -  

Universität Dortmund




## Improving predictability for caches

- Loop caches
- Mapping code to less used part(s) of the index space
- Cache locking/freezing
- 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”

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 99 -  

Universität Dortmund

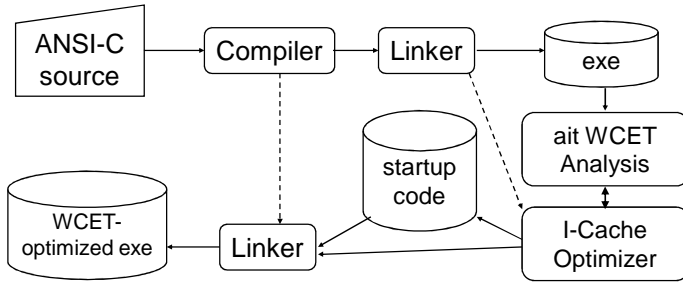
## Locking of I-Caches

Many caches allow “locking” or “freezing” (no replacements).  
Can be used to improve timing predictability:

- Load promising Functions into Cache




Optimizations:

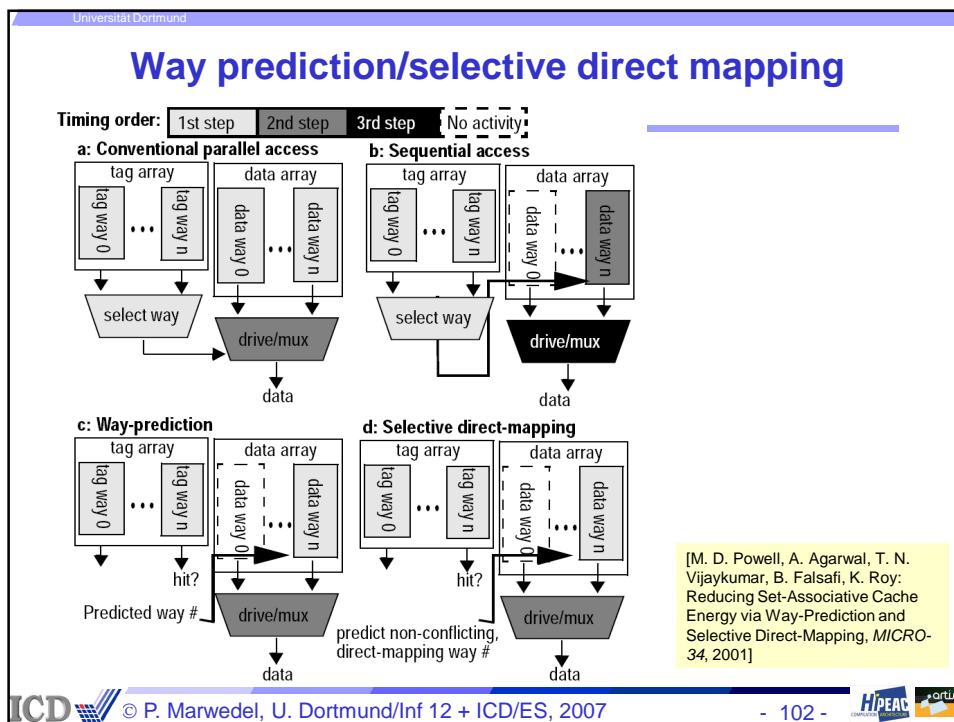
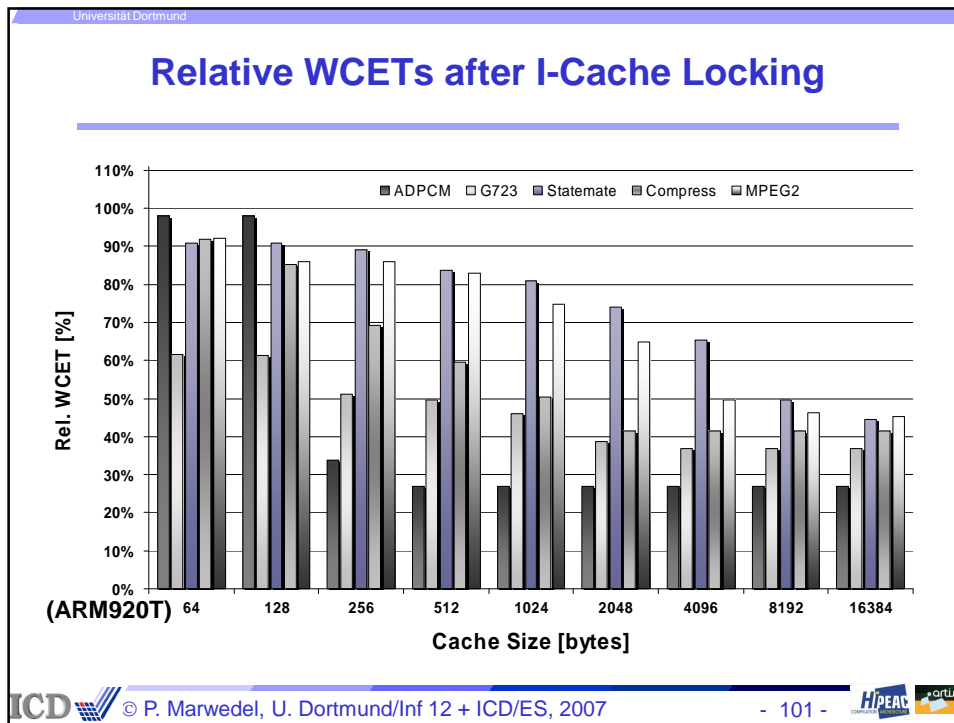
- Worst case paths can change during optimization
- Requires coupling of timing analysis tool and compiler




```
graph LR
    Source[ANSI-C source] --> Compiler[Compiler]
    Compiler --> Linker1[Linker]
    Linker1 --> Exe[(exe)]
    Exe --> WCET[WCET-optimized exe]
    Exe --> WCET_Analysis[WCET Analysis]
    WCET_Analysis <--> IOpt[I-Cache Optimizer]
    IOpt --> Linker2[Linker]
    Linker2 --> Startup[(startup code)]
    Linker2 --> Linker1
    Linker2 --> WCET_Analysis
```

[Heiko Falk, Sascha Plazar, Henrik Theiling:  
Compile-Time Decided  
Instruction Cache  
Locking Using Worst-  
Case Execution Paths,  
CODES/ISSS, 2007]

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 100 -  



## Outline

- Introduction, Motivation
- Scratch-pad memories
  - Characteristics
  - Non-overlaying allocation, partitions, stack, heap, WCET
  - Overlaying (dynamic) allocation
  - Multiple processes & processors
- Caches 
  - Predictability, loop transformations
  - Prefetching, layout transformations
- Other levels of the memory hierarchy
  - Array folding
  - ADL-driven compilation for hierarchies
  - Rewriting algorithms for memory hierarchies
- Conclusion, Additional references

## Locality of Reference

### Spatial Locality:

- The concept that likelihood of referencing an item is higher if an item close to it was just referenced.

### Temporal Locality:

- The concept that a item that is referenced at one point in time will be referenced again in the near future.

```
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];  
    }  
  }  
}
```

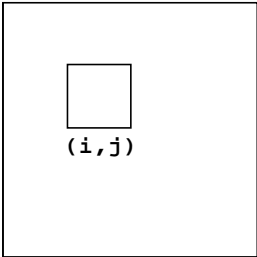
Universität Dortmund

## Locality: Matrix Multiplication

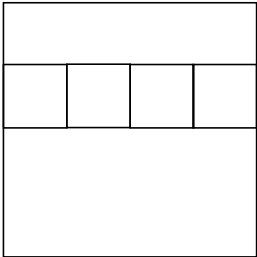
---

```
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]; }}}
```

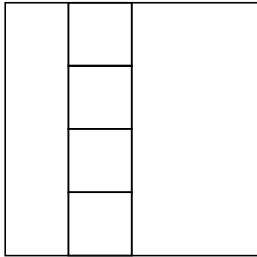
$I_1 = (i, j, k)$  and  $I_2 = (i, j, k+1)$   
 $c[i][j] \rightarrow$  Temporal Locality



**c**






**a**



**b**

\*

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 105 -  

Universität Dortmund

## Row major order vs. column major order allocation

---

Row major order (C)

Array p[j][k]

Column major order (FORTRAN)

j=0 {

...
k=0
k=1
...

j=1 {

k=0
k=1
...

j=2 {

k=0
k=1
...

...

k=0 {

...
j=0
j=1
...




k=1 {

j=0
j=1
...

k=2 {

j=0
j=1
...

...

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 106 -  

Universitat Dortmund

## Locality: Matrix Multiplication (2)

---

$I_1 = (i, j, k)$  and  $I_2 = (i, j, k+1)$   
 $c[i][j] \rightarrow$  Temporal Locality (layout independent)  
 $a[i][k] \rightarrow$  Spatial Locality (Row major order)  
 $b[k][j] \rightarrow$  Spatial Locality (Column major order)

**c**

**a**

**b**

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 107 -

Universitat Dortmund

## Data Layout Transformations - Example

---

Change layout of arrays according to the access pattern

```
int c[4][4];
for (i=0; i<4; i++) {
  for (j=0; j<2; j++) {
    c[i][j]= ...; }}
```

⇒

```
int c[2][8];
for (i=0; i<4; i++) {
  for (j=0; j<2; j++) {
    c[2*i+j]= ...; }}
```

Poor cache behavior

Good cache behavior

Increases the spatial locality of array accesses

[V. Loechner et al.: Precise Data Locality Optimization of Nested Loops, *The Journal of Supercomputing*, Vol 21, 2002]

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 108 -

## Cache Aware Compiler Optimizations

Performance of a cache depends upon locality of references.

Data Cache:

- Loop Restructuring Transformations → Temporal Locality
- Data Layout Transformations → Spatial Locality

Instruction Cache:

- Code Layout Transformations → Spatial Locality

## Loop Restructuring Transformations - Loop Interchange -

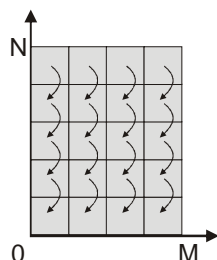
Assumption: row major order

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

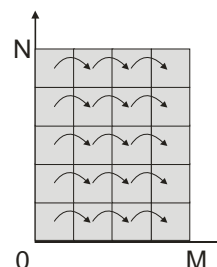


```
for (j=0; j<N; j++) {  
  for (i=0; i<M; i++) {  
    c[j][i] += ...  
  }  
}
```

Poor cache behavior



Good cache behavior



Best performance if innermost loop  
corresponds to rightmost array index

[D. F. Bacon et al.: Compiler Transformations for High-  
Performance Computing, ACM Computing Surveys, 1994]

Universität Dortmund

## Loop interchange: Example

```

...#define iter 400000
int a[20][20][20];
void computeijk() {int i,j,k;
    for (i = 0; i < 20; i++) {
        for (j = 0; j < 20; j++) {
            for (k = 0; k < 20; k++) {
                a[i][j][k] += a[i][j][k];}}}
void computeikj() {int i,j,k;
    for (i = 0; i < 20; i++) {
        for (j = 0; j < 20; j++) {
            for (k = 0; k < 20; k++) {
                a[i][k][j] += a[i][k][j];}}}...
start=time(&start);for(z=0;z<iter;z++)computeijk();
end=time(&end);
printf("ijk=%16.9f\n",1.0*difftime(end,start));

```

(SUIF interchanges array indexes instead of loops; low resolution time calls)

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 111 - H/PEAC

Universität Dortmund

## Results:

strong influence of the memory architecture

Loop structure: i j k Dramatic impact of locality

Processor	Ti C6xx	Sun SPARC	Intel Pentium
reduction to [%]	~ 57%	35%	3.2 %

[Till Buchwald, Diploma thesis, Univ. Dortmund, Informatik 12, 12/2004]

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 112 - H/PEAC



Universität Dortmund

## Loop tiling/loop blocking Original version

```

for (i=1; i<=N; i++)
  for(k=1; k<=N; k++){
    r=X[i,k]; /* to be allocated to a register*/
    for (j=1; j<=N; j++)
      Z[i,j] += r* Y[k,j]
  } % Never reusing information in the cache for Y and Z if N
    is large or cache is small (O(N³) references for Z and Y).
  
```

© P.Marwedel, U. Dortmund, Informatik 12, 2006 - 113 -

Universität Dortmund

## Loop tiling/loop blocking tiled version

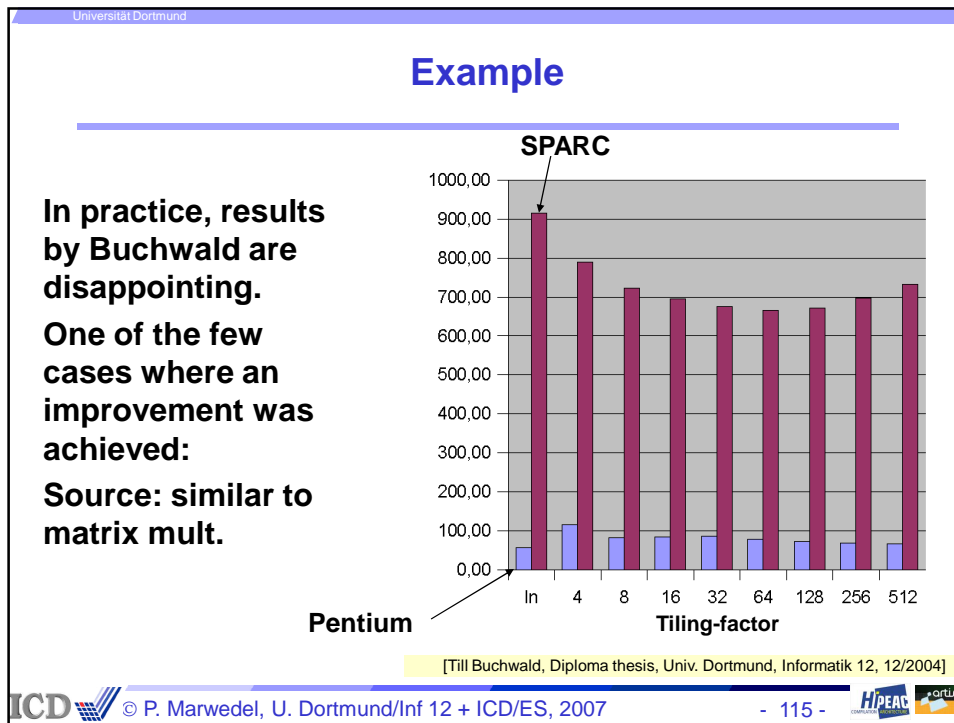
```

for (kk=1; kk<= N; kk+=B)
  for (jj=1; jj<= N; jj+=B)
    for (i=1; i<= N; i++)
      for (k=kk; k<= min(kk+B-1,N); k++){
        r=X[i][k]; /* to be allocated to a register*/
        for (j=jj; j<= min(jj+B-1, N); j++)
          Z[i][j] += r* Y[k][j]
      }
  
```

Reuse factor of B  
for Z and  
N for Y,  
  
O(N³/B) accesses  
to main memory  
for Z

Same elements for  
next iteration of i

© P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 114 -



Universität Dortmund

## Loop fusion (merging), loop fission

```

for(j=0; j<=n; j++)
  p[j]= ... ;
for (j=0; j<=n; j++) ,
  p[j]= p[j] + ...
        
```

```

for (j=0; j<=n; j++)
  {p[j]= ... ;
   p[j]= p[j] + ...}
        
```

Loops small enough to allow zero overhead

Loops

**Better locality** for access to p.

Better chances for parallel execution.

Which of the two versions is best?

Architecture-aware compiler should select best version.

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 116 - H/PEAC partit

## Example: loops with reuse

```
#define size 30
#define iter 40000
int a[size][size];
float b[size][size];
```

```
void ss2() {int i,j;
  for(i=0;i<size;i++) {
    for(j=0;j<size;j++) {
      a[i][j] += 17;}}
  for(i=0;i<size;i++){
    for(j=0;j<size;j++) {
      b[i][j]-=a[i][j];}}}
```

```
void ms2() {int i,j;
  for(i=0;i<size;i++) {
    for(j=0;j<size;j++) {
      a[i][j]+= 17;}
    for(j=0;j<size;j++){
      b[i][j]-=a[i][j];}}}
```

```
void mm2() {int i,j;
  for(i=0;i<size;i++) {
    for(j=0;j<size;j++) {
      a[i][j]+=17;
      b[i][j]-=a[i][j];}}}
```

## Loop unrolling

```
for (j=0; j<=n; j++)
  p[j]= ... ;
```



```
for (j=0; j<=n; j+=2)
  {p[j]= ... ; p[j+1]= ... }
```

factor = 2

**Better locality** for access to p.

Less branches per execution of the loop. More opportunities for optimizations.

Tradeoff between code size and improvement.

Extreme case: completely unrolled loop (no branch).

## Example: matrixmult

```
#define s 30
#define iter 4000
int
a[s][s],b[s][s],c[s]
[s];
void compute(){int
i,j,k;
for(i=0;i<s;i++){
for(j=0;j<s;j++){
for(k=0;k<s;k++){
c[i][k]+=
a[i][j]*b[j][k];
}}}}

extern void compute2()
{int i, j, k;
for (i = 0; i < 30; i++) {
for (j = 0; j < 30; j++) {
for (k = 0; k <= 28; k += 2)
{{int *suif_tmp;
suif_tmp = &c[i][k];
*suiif_tmp=
*suiif_tmp+a[i][j]*b[j][k];}
{int *suif_tmp;
suif_tmp=&c[i][k+1];
*suiif_tmp=*suiif_tmp
+a[i][j]*b[j][k+1];
}}}}
return;}
```

## Summary Lecture 3

- Overlaying SPM allocation
  - Multiple processes
    - Dynamic set of processes
    - MMU-based approach (SNU)
  - Multiple processors
    - Verma's approach ( $\Phi$ MPARM)
    - Kandemir's MPSoC approach
- Caches
  - Locked caches
  - Locality of reference
  - Loop transformations, tiling



## Compilation for Caches

# Compilation for other memory hierarchy levels

## Outline

- Introduction, Motivation
- Scratch-pad memories
  - Characteristics
  - Non-overlaying allocation, partitions, stack, heap, WCET
  - Overlaying (dynamic) allocation
  - Multiple processes & processors
- Caches
  - Predictability, loop transformations
  - ➡ - Prefetching, layout transformations
- Other levels of the memory hierarchy
  - Array folding
  - ADL-driven compilation for hierarchies
  - Rewriting algorithms for memory hierarchies
- Conclusion, Additional references




Universität Dortmund

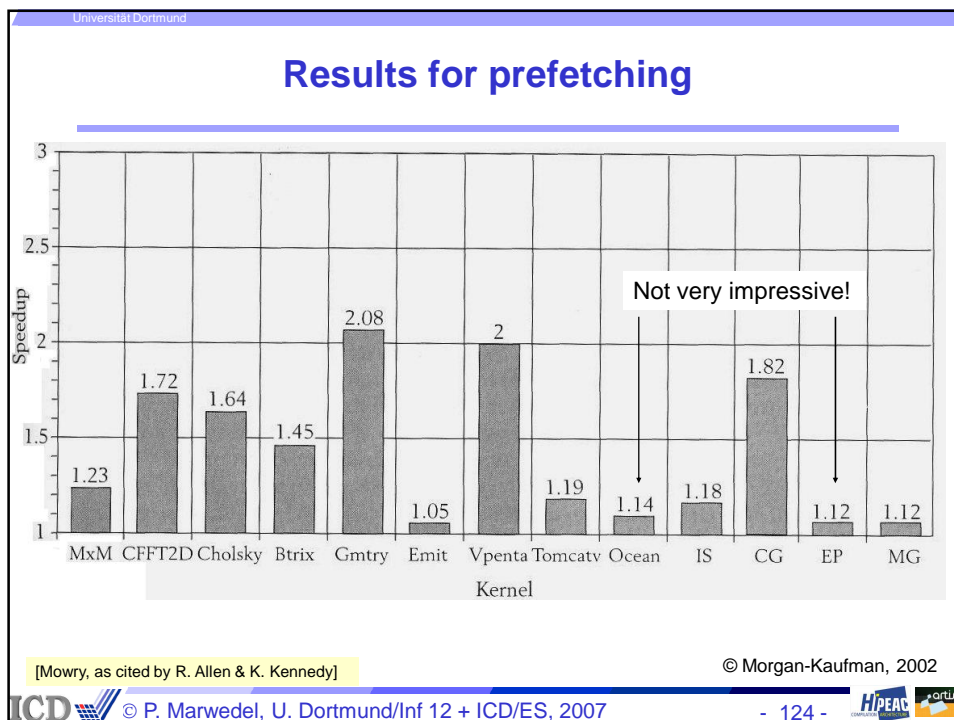
## Prefetching

---

- Prefetch instructions load values into the cache  
Pipeline not stalled for prefetching
- Prefetching instructions introduced in ~1985-1995
- Potentially, all miss latencies can be avoided
- Disadvantages:
  - Increased # of instructions
  - Potential premature eviction of cache line
  - Potentially pre-loads lines that are never used
- Steps
  - Determination of references requiring prefetches
  - Insertion of prefetches (early enough!)

[R. Allen, K. Kennedy: Optimizing Compilers for Modern Architectures, Morgan-Kaufman, 2002]

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 123 -  



Universität Dortmund

## Optimization for exploiting processor-memory interface: Problem Definition (1)

XScale is stalled for 30% of time, but each stall duration is small

- Average stall duration = 4 cycles
- Longest stall duration < 100 cycles

Break-even stall duration for profitable switching

- 360 cycles

Maximum processor stall

- < 100 cycles

NOT possible to switch the processor to IDLE mode

**Processor Stall Durations**

[A. Shrivastava, E. Earlie, N. Dutt, A. Nicolau: Aggregating processor free time for energy reduction, Intern. Conf. on Hardware/Software Codesign and System Synthesis (CODES/ISSS), 2005, pp. 154-159] Based on slide by A. Shrivastava

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 125 - H/PEAC

Universität Dortmund

## Optimization for exploiting processor-memory interface: Problem Definition (2)

- CT (Computation Time): Time to execute an iteration of the loop, assuming all data is present in the cache
- DT (Data Transfer Time): Time to transfer data required by an iteration of a loop between cache and memory

Consider the execution of a memory-bound loop (DT > CT)

- Processor has to stall

```
for (int i=0; i<1000; i++)
  c[i] = a[i] + b[i];
```

Activity

Time

Processor Activity

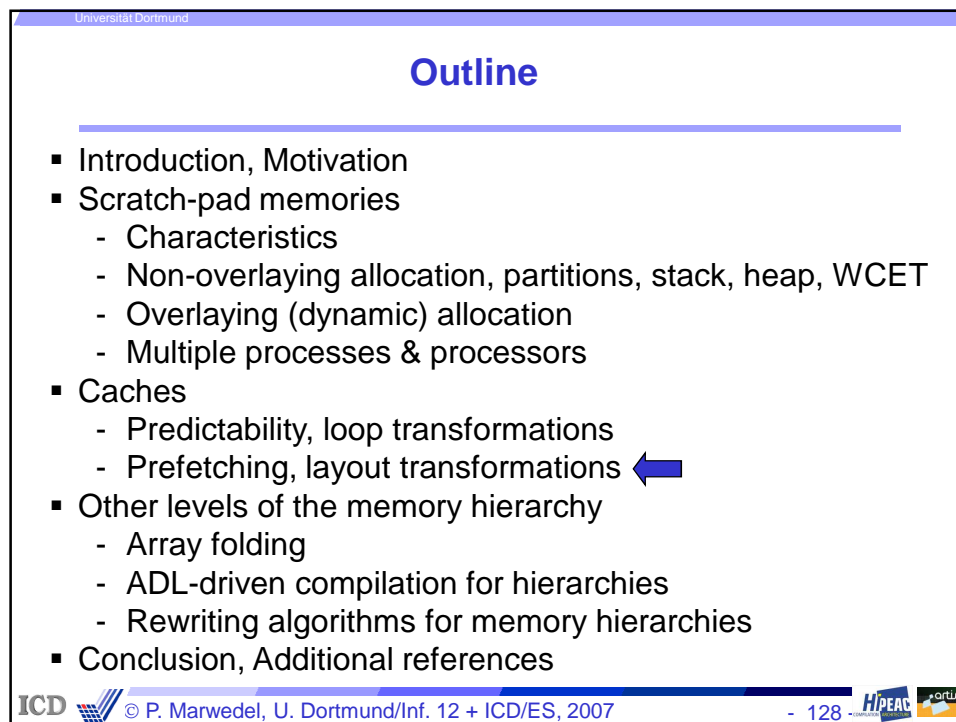
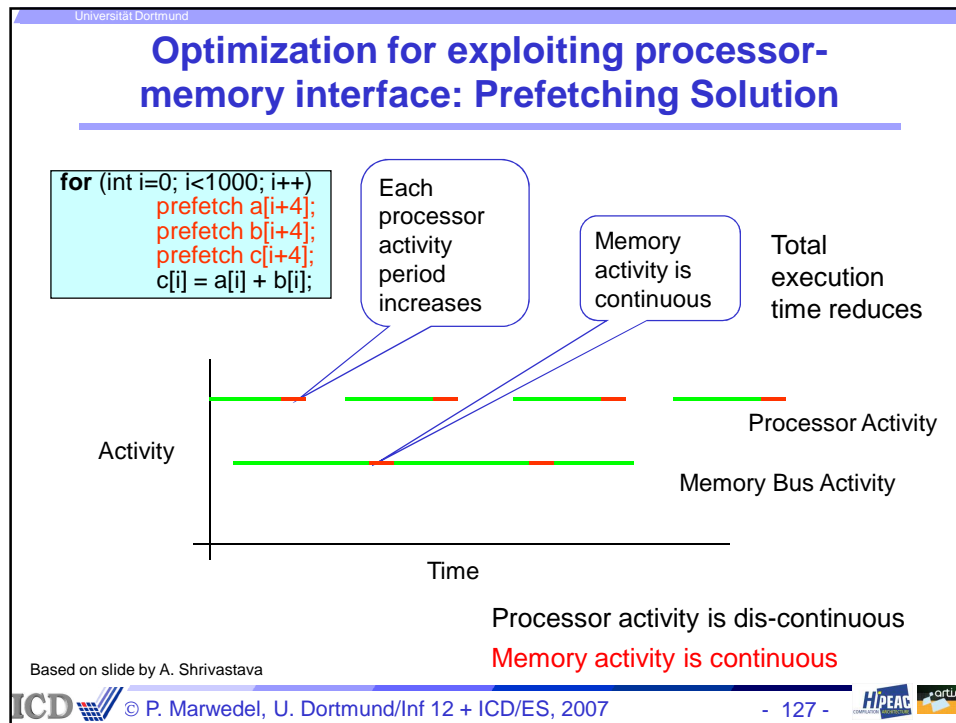
Memory Bus Activity

Processor activity is dis-continuous

Memory activity is dis-continuous

Based on slide by A. Shrivastava

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 126 - H/PEAC





Universitat Dortmund




## Code Layout Transformations (1)

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

(1100)	$f_1$
(900)	$f_2$
(400)	$f_3$
(2000)	$f_4$
(700)	$f_5$

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

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 129 -  

Universitat Dortmund

## 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

```

graph TD
    f4((f4)) --> f2((f2))
    f4 --> f5((f5))
    f2 --> f1((f1))
    f2 --> f3((f3))
  
```

(2000)	$f_4$
(1100)	$f_1$
(900)	$f_2$
(700)	$f_5$
(400)	$f_3$

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 130 -  

Universitat Dortmund

## Code Layout Transformations (3)

Call-Graph Based Algorithm: (2000)

- 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.

```

graph TD
    f4((f4)) --> f2((f2))
    f4 --> f5((f5))
    f2 --> f1((f1))
    f2 --> f3((f3))
  
```

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

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 131 - H/PEAC artur

Universitat Dortmund

## Code Layout Transformations (3)

Call-Graph Based Algorithm: (2000)

- Create weighted call-graph.
- Place functions according to weighted depth-first traversal. (900)

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

Increases spatial locality.

```

graph TD
    f4((f4)) --> f2((f2))
    f4 --> f5((f5))
    f2 --> f1((f1))
    f2 --> f3((f3))
  
```

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 132 - H/PEAC artur

Universitat Dortmund

## 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)

(900)

(1100)

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 133 -

Universitat Dortmund

## 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)

(900)

(1100)

(400)

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 134 -

Universitat Dortmund

## 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.

```

graph TD
    f4((f4)) --> f2((f2))
    f4 --> f5((f5))
    f2 --> f1((f1))
    f2 --> f3((f3))
    
```

(2000) f<sub>4</sub>

(900) f<sub>2</sub>

(1100) f<sub>1</sub>

(400) f<sub>3</sub>

(700) f<sub>5</sub>

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 135 -

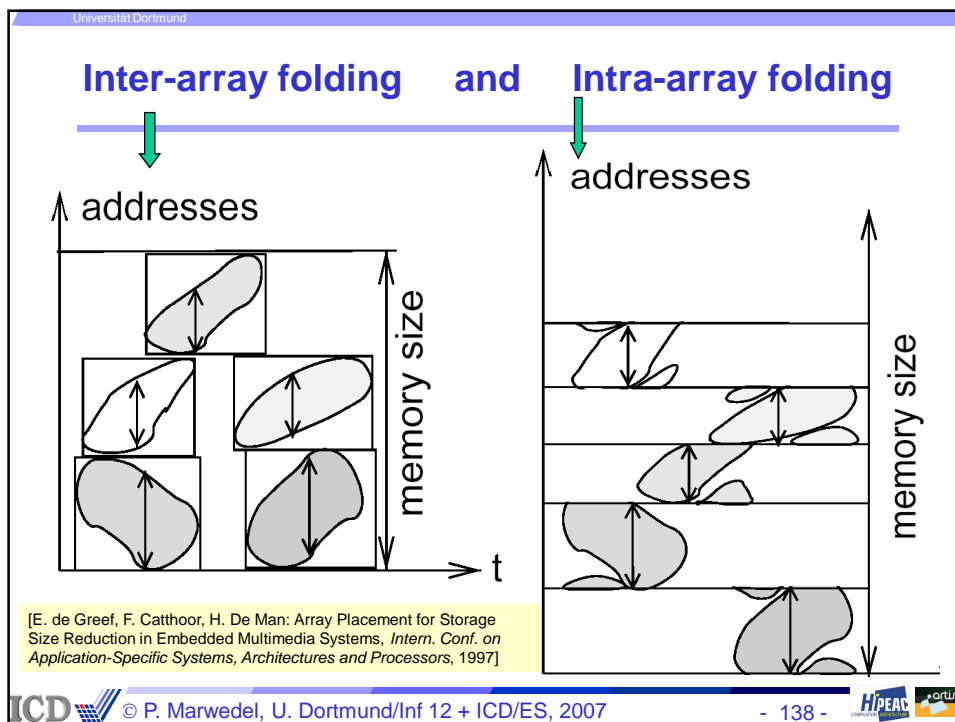
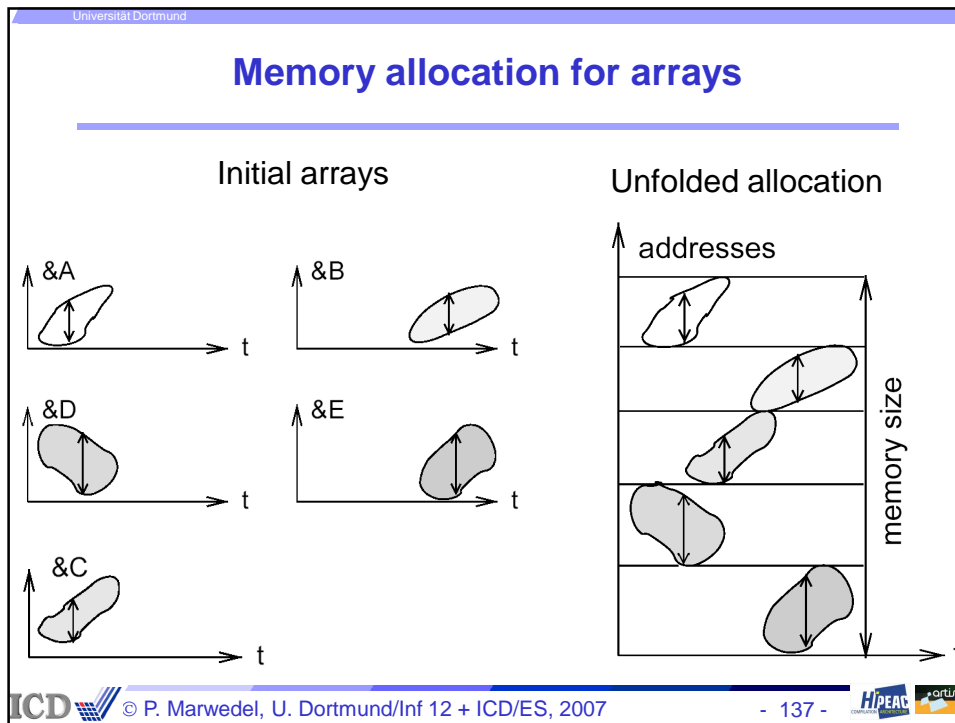
Universitat Dortmund

## Outline

---

- Introduction, Motivation
- Scratch-pad memories
  - Characteristics
  - Non-overlying allocation, partitions, stack, heap, WCET
  - Overlaying (dynamic) allocation
  - Multiple processes & processors
- Caches
  - Predictability, loop transformations
  - Prefetching, layout transformations
- Other levels of the memory hierarchy
  - ➡ - Array folding
  - ADL-driven compilation for hierarchies
  - Rewriting algorithms for memory hierarchies
- Conclusion, Additional references

ICD © P. Marwedel, U. Dortmund/Inf. 12 + ICD/ES, 2007 - 136 -



Universität Dortmund

## Requires optimization of address computations

- Array folding implemented in IMEC's DTSE optimization.
- Leads to costly div and mod ops.
- ADOPT address optimizations remove these operations  
E.g.:  $\text{mod} \rightarrow ++$  and reset on pointers (indexes)

```

for(i=0; i<20; i++)
  B[i%4];
  
```

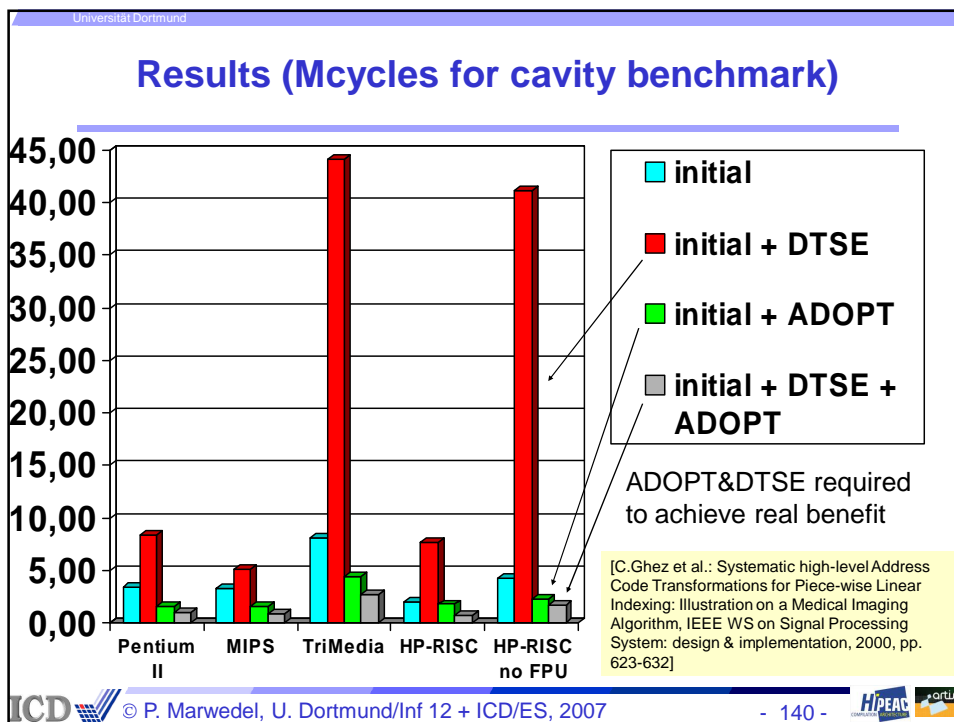
↓

```

tmp=0;
for(i=0; i<20; i++)
{
  if(tmp>=4)
    tmp-=4;
  B[tmp];
  tmp++;
}
  
```

*Architecture-aware compiler should find best transformation*

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 139 - H/PEAC artur






Universität Dortmund

## Outline

---

- Introduction, Motivation
- Scratch-pad memories
  - Characteristics
  - Non-overlying allocation, partitions, stack, heap, WCET
  - Overlaying (dynamic) allocation
  - Multiple processes & processors
- Caches
  - Predictability, loop transformations
  - Prefetching, layout transformations
- Other levels of the memory hierarchy
  - Array folding
  - ➔ - ADL-driven compilation for hierarchies
  - Rewriting algorithms for memory hierarchies
- Conclusion, Additional references

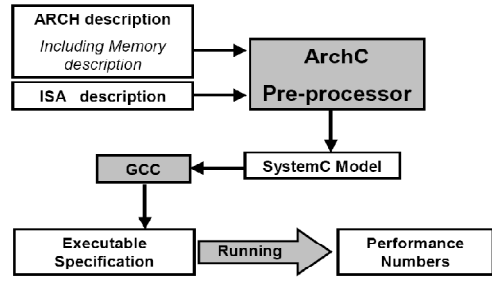
ICD  © P. Marwedel, U. Dortmund/Inf. 12 + ICD/ES, 2007 - 141  

Universität Dortmund

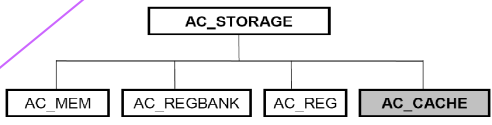
## Memory hierarchy description languages: ArchC

---




Consists of description of ISA and HW architecture  
Extension of SystemC (can be generated from ArchC):



Storage class structure

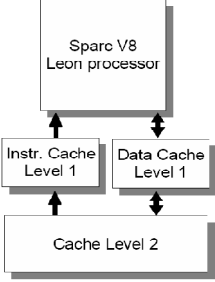


[P. Viana, E. Barros, S. Rigo, R. Azevedo, G. Araújo:  
Exploring Memory Hierarchy with ArchC, 15th  
Symposium on Computer Architecture and High  
Performance Computing, 2003, pp. 2 – 9]

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 142 -  

Universität Dortmund

## Example: Description of a simple cache-based architecture



```

AC_ARCH(leon) {

    ac_cache  icache("dm", 128, "wt")
    ac_cache  dcache("2w", 64, 4, "wt", "lru")
    ac_cache  ul2cache("dm", 4k, "wt")




    ac_regbank RB:520;
    ac_reg PRS, Y, WIM;

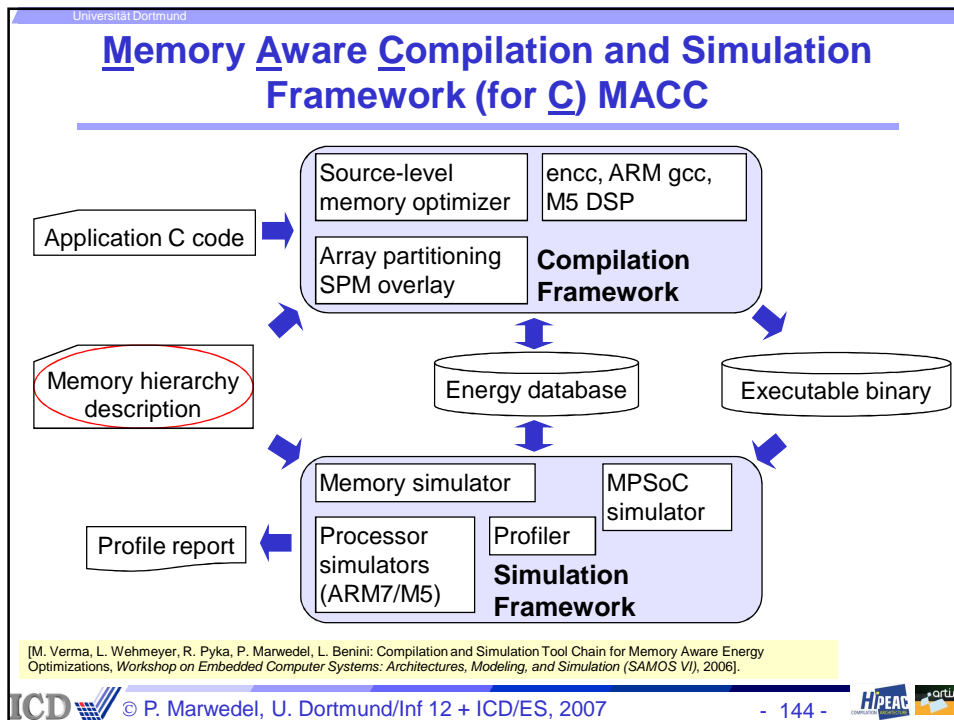
    ac_pipe   pipe = {IF, ID, EX, MEM, WB};

    ARCH_CTOR(leon) {
        ac_isa("leon_isa.ac");

        icache.bindTo( ul2cache ); //Memory hierarchy
        dcache.bindTo( ul2cache ); //construction
    };
};

```

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 143 -  





Universität Dortmund

## Memory architecture description @ MACCv2

- Query can include address, time stamp, value, ...
- Query can request energy, delay, stored values
- Query processed along a chain of HW components, incl. busses, ports, address translations etc., each adding delay & energy
- API query to model simplifies integration into compiler
- External XML representation

REQ: Energy= ?, Cycles= ?

CPU1 → ASPC-1 (IFETCH, DRD, DWR, MAINAS) → ASPC-B (0...3fff) → ASPC-M (0...ffff) → MM

Energy costs: +1 → Energy, +1 → Energy, +10 → Energy  
Cycle costs: +0 → Cycles, +2 → Cycles, +5 → Cycles

[R. Pyka et al.: Versatile System level Memory Description Approach for embedded MPSoCs, University of Dortmund, Informatik 12, 2007]

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 145 - HPEAC artur

Universität Dortmund

## Controlling tool chain generation through an architecture description language (ADL): EXPRESSION

Overall information flow

Architecture Specification (English Document) → ADL Specification

ADL Specification branches into:

- Generic VHDL Models → HDL Generator → Hardware Model → Synthesis
- Compiler Generator → Compiler (receiving Application Programs) → Binary
- Generic C++ Models → Simulator Generator → Simulator

Feedback loops:

- Feedback (Area, Power, Clock Frequency) from Synthesis to Architecture Specification
- Feedback (Performance, Code Size) from Simulator to Architecture Specification

Legend: — Automatic, - - - Manual, ····· Feedback

[P. Mishra, A. Shrivastava, N. Dutt: Architecture description language (ADL)-driven software toolkit generation for architectural exploration of programmable SOCs, ACM Trans. Des. Autom. Electron. Syst. (TODAES), 2006, pp. 626-658]

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 146 - HPEAC artur

Universität Dortmund

## Description of Memories in EXPRESSION

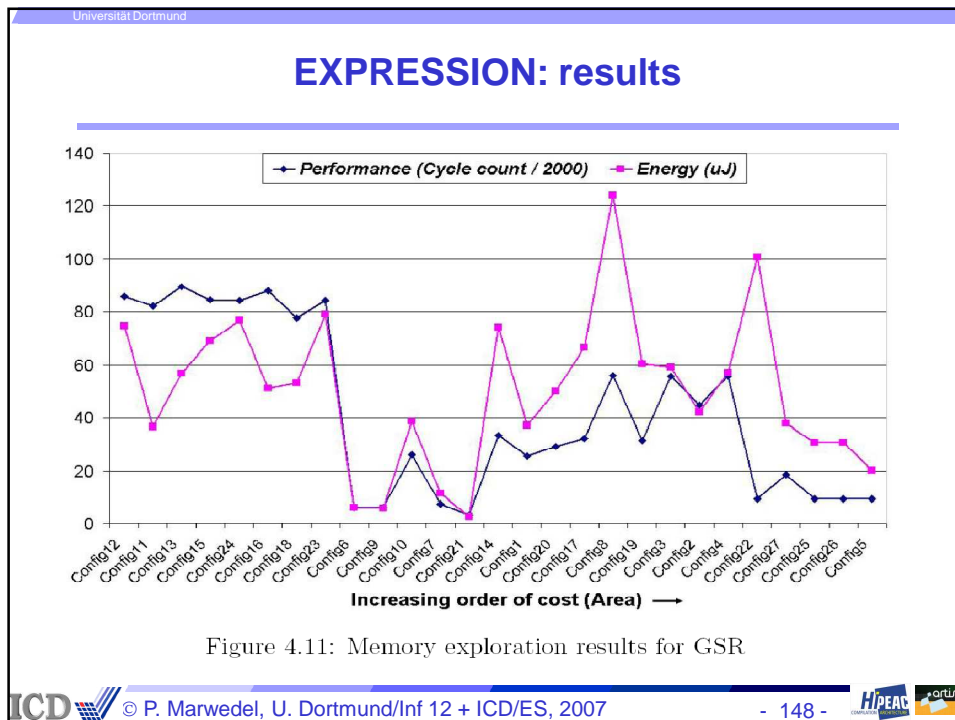
Generic approach, based on the analysis of a wide range of systems;

Used for verification.

```
(STORAGE_SECTION
(DataL1
(TYPE DCACHE) (WORDSIZE 64)
(LINESIZE 8) (NUM_LINES 1024)
(ASSOCIATIVITY 2) (READ_LATENCY 1) ...
(REPLACEMENT_POLICY LRU)
(WRITE_POLICY WRITE_BACK)
)
(ScratchPad
(TYPE SRAM) (ADDRESS_RANGE 0 4095) ....
)
(SB
(TYPE STREAM_BUFFER) .....
)
(InstL1
(TYPE ICACHE) .....
)
(L2
(TYPE DCACHE) .....
)
(MainMemory
(TYPE DRAM)
)
)
(Connect
(TYPE CONNECTIVITY)
(CONNECTIONS
(InstL1, L2) (DataL1, SB) (SB, L2)
(L2, MainMemory)
)
)
))
```

© P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007

- 147 -








Universität Dortmund

## Outline

---

- Introduction, Motivation
- Scratch-pad memories
  - Characteristics
  - Non-overlying allocation, partitions, stack, heap, WCET
  - Overlaying (dynamic) allocation
  - Multiple processes & processors
- Caches
  - Predictability, loop transformations
  - Prefetching, layout transformations
- Other levels of the memory hierarchy
  - Array folding
  - ADL-driven compilation for hierarchies
- ➔ - Rewriting algorithms for memory hierarchies
- Conclusion, Additional references

ICD  © P. Marwedel, U. Dortmund/Inf. 12 + ICD/ES, 2007 - 151  

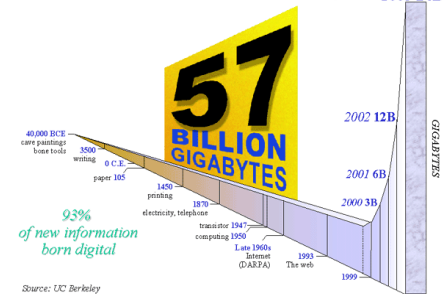
Universität Dortmund

## Memory hierarchies beyond main memory

---

- Massive datasets are being collected everywhere
- Storage management software is billion-\$ industry

More New Information Over Next 2 Years Than in All Previous History






Source: UC Berkeley  
EMC Copyright 2001

Examples (2002):

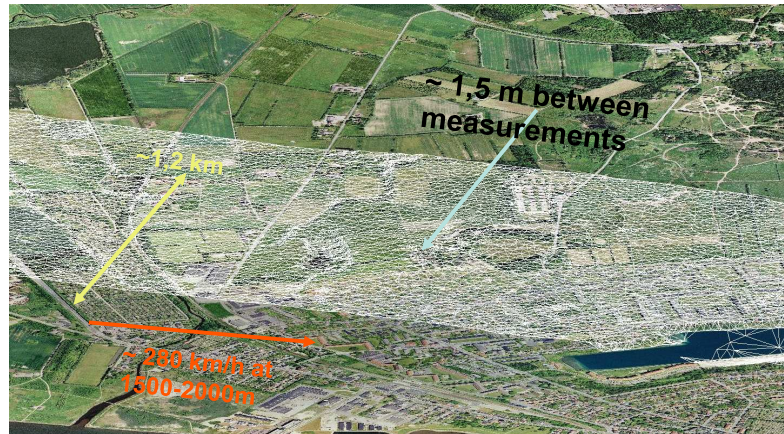
- **Phone**: AT&T 20TB phone call database, wireless tracking
- **Consumer**: WalMart 70TB database, buying patterns
- **WEB**: Web crawl of 200M pages and 2000M links, Akamai stores 7 billion clicks per day
- **Geography**: NASA satellites generate 1.2TB per day

© Lars Arge, I/O-Algorithms, <http://www.daimi.au.dk/~large/ioS07/>

ICD  © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 152 -  

## Example: LIDAR Terrain Data

- COWI A/S (and others) is currently scanning Denmark

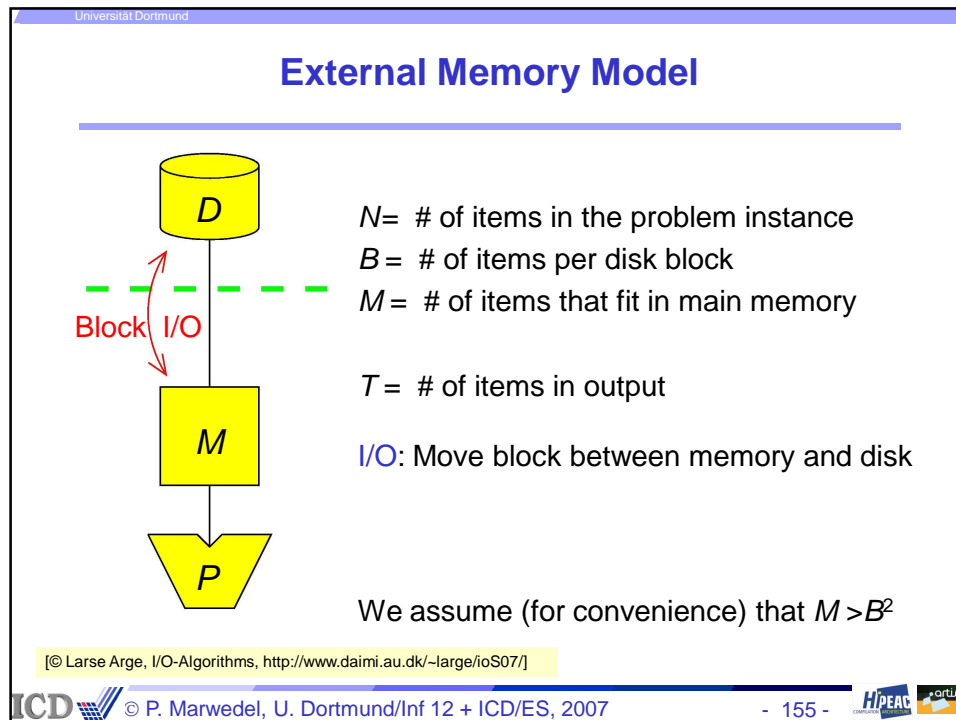


[© Larse Arge, I/O-Algorithms, <http://www.daimi.au.dk/~larse/ioS07/>]

## Application Example: Flooding Prediction



[© Larse Arge, I/O-Algorithms, <http://www.daimi.au.dk/~larse/ioS07/>]



Universität Dortmund

## Scalability Problems: Block Access Matters

---

- **Example:** Reading an array from disk
  - Array size  $N = 10$  elements
  - Disk block size  $B = 2$  elements
  - Main memory size  $M = 4$  elements (2 blocks)

1	5	2	6	3	8	9	4	7	10
---	---	---	---	---	---	---	---	---	----

Algorithm 1:  $N=10$  I/Os

1	2	10	9	5	6	3	4	8	7
---	---	----	---	---	---	---	---	---	---

Algorithm 2:  $N/B=5$  I/Os

- Difference between  $N$  and  $N/B$  large since block size is large
  - **Example:**  $N = 256 \times 10^6$ ,  $B = 8000$ ,  $1ms$  disk access time
    - $\Rightarrow N$  I/Os take  $256 \times 10^3 \text{ sec} = 4266 \text{ min} = 71 \text{ hr}$
    - $\Rightarrow N/B$  I/Os take  $256/8 \text{ sec} = 32 \text{ sec}$

© Larse Arge, I/O-Algorithms, <http://www.daimi.au.dk/~larse/ioS07/>

ICD © P. Marwedel, U. Dortmund/Inf 12 + ICD/ES, 2007 - 156 -

## Re-writing algorithms for memory hierarchies

Analysis of algorithm complexity mostly using the *RAM* (random access machine; const. mem. acc. times) model outdated

☞ take memory hierarchies explicitly into account.

Example:

- Usually, divide-&-conquer algorithms are good.
- “Cache”-oblivious algorithms (are good for any size of the faster memory and any block size). Assuming
  - Optimal replacement (Belady’s algorithm)
  - 2 Memory levels considered (there can be more)
  - Full associativity
  - Automatic replacement

[Piyush Kumar: Cache Oblivious Algorithms, in: U. Meyer et al. (eds.): Algorithms for Memory Hierarchies, *Lecture Notes in Computer Science*, Volume 2625, 2003, pp. 193-212]

[Naila Rahman: Algorithms for Hardware Caches and TLB, in: U. Meyer et al. (eds.): Algorithms for Memory Hierarchies, *Lecture Notes in Computer Science*, Volume 2625, 2003, pp. 171-192]

Unlikely to be ever automatic

## Fundamental Bounds

### Internal

- **Scanning:**  $N$
- **Sorting:**  $N \log N$
- **Permuting**  $N$
- **Searching:**  $\log_2 N$

### Note:

- Linear I/O:  $O(N/B)$
- Permuting not linear
- Permuting and sorting bounds are equal in all practical cases
- $B$  factor VERY important:  $\frac{N}{B} < \frac{N}{B} \log_{M/B} \frac{N}{B} \ll N$

### External

$$\frac{N}{B}$$

$$\frac{N}{B} \log_{M/B} \frac{N}{B}$$

$$\min\{N, \frac{N}{B} \log_{M/B} \frac{N}{B}\}$$

$$\log_B N$$

[© Larse Arge, I/O-Algorithms, <http://www.daimi.au.dk/~larse/ioS07/>]



## Additional References (1)

- D. Page: Effective Use of Partitioned Cache Memories, *PhD Thesis, Department of Computer Science, University of Bristol*, UK 2001.
- O. S. Unsal et al. Cool-Cache: A Compiler-Enabled Energy Efficient Data Caching Framework for Embedded/ Multimedia Processors, *ACM Transactions on Embedded Computing Systems*, Vol. 2, 2003
- J. Kin et al.: The Filter Cache: An Energy Efficient Memory Structure, *30th Annual International Symposium on Microarchitecture (MICRO)*, 1997.
- L. H. Lee et al. Instruction Fetch Energy Reduction Using Loop Caches For Embedded Applications with small Tight Loops, *International Symposium on Low Power Electronics and Design (ISLPED)*, 1999.
- S. C. A. Gordon-Ross, F. Vahid: Exploiting Fixed Programs in Embedded Systems: A Loop Cache Example, *Computer Architecture Letters*, 1, 2002.
- S. Steinke: Untersuchung des Energieeinsparungspotenzials in eingebetteten Systemen durch energieoptimierende Compilertechnik (in German), *PhD Thesis, University of Dortmund, Germany*, 2003.

## Additional References (2)

- P. R. Panda, N. Dutt, A. Nicolau: Memory Issues in Embedded Systems-On-Chip, *Kluwer Academic Publishers, Norwell, MA*, 1999.
- S. Steinke et al.: Assigning Program and Data Objects to Scratchpad for Energy Reduction, *Design Automation and Test in Europe Conference (DATE'02)*, 2002.
- L. Wehmeyer, P. Marwedel: Influence of On-chip Scratchpad Memories on WCET, *4th International Workshop on Worst-Case Execution Time Analysis (WCET)*, 2004
- L. Wehmeyer, P. Marwedel: Influence of Memory Hierarchies on Predictability for Time Constrained Embedded Software, *Design Automation and Test in Europe Conference (DATE'05)*, 2005
- L. Wehmeyer, U. Helmig, P. Marwedel: Compiler-optimized Usage of Partitioned Memories, *3rd Workshop on Memory Performance Issues (WMPI)*, 2004.
- F. Angiolini, L. Benini, A. Caprara: Polynomial-Time Algorithm for On-Chip Scratchpad Memory Partitioning, *International Conference on Compilers, Architecture and Synthesis of Embedded Systems (CASES'03)*, 2003

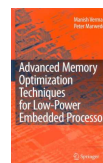


## Additional References (3)

- J. Chame Shin, M. W. Hall: A Compiler Algorithm for Exploiting Page-Mode Memory Accesses in Embedded-DRAM Devices, *Workshop on Media and Stream Processors*, 2002
- M. Verma, P. Marwedel: Overlay of Scratchpad Memory for Low Power Embedded Processors, *IEEE Transactions on Very Large Scale Integration Systems (TVLSI)*, Vol.14, 2006.
- E. Brockmeyer et al.: Layer Assignment Techniques for Low Energy in Multi-Layered Memory Organisations, *Design, Automation and Test in Europe Conference (DATE)*, 2003.
- M. Verma, K. Petzold, L. Wehmeyer, H. Falk, P. Marwedel: Scratchpad Sharing Strategies for Multiprocess Embedded Systems: A First Approach, *IEEE 3rd Workshop on Embedded Real-Time Multimedia (ESTIMedia)*, 2005.
- ❖ L. Wehmeyer, P. Marwedel: Fast, Efficient and Predictable Memory Accesses Optimization Algorithms for Memory Architecture Aware Compilation, *Springer*, 2006
- ❖ M. Verma, Peter Marwedel: Advanced Memory Optimization Techniques for Low-Power Embedded Processors, *Springer*, 2007

## New Books

- Lars Wehmeyer, Peter Marwedel: Fast, Efficient and Predictable Memory Accesses, *Springer*, 2006
- Manish Verma, Peter Marwedel: Advanced Memory Optimization Techniques for Low-Power Embedded Processors, *Springer*, May 2007
- Peter Marwedel: Embedded System Design, *Springer*, 2005 (German edition: 2007)



## Summary lecture 4

- Caches
  - Prefetching
  - Code layout transformations
- Other memory hierarchy levels
  - Array folding
  - ADL-driven compilation for memory hierarchies
  - Rewriting algorithms for memory hierarchies

## Overall summary

- Introduction, Motivation
- Scratch-pad memories
  - Characteristics
  - Non-overlaying allocation, partitions, stack, heap, WCET
  - Overlaying (dynamic) allocation
  - Multiple processes & processors
- Caches
  - Predictability, loop transformations
  - Prefetching, layout transformations
- Other levels of the memory hierarchy
  - Array folding
  - ADL-driven compilation for hierarchies
  - Rewriting algorithms for memory hierarchies
- Conclusion, Additional references

Credit goes to all students involved in the design of SPM tools at Dortmund, in particular to Stefan Steinke, Robert Pyka, Manish Verma and Lars Wehmeyer.

**Thank you for your attention!**  
**Questions?**