

PICO: ASIC Synthesis from C

Rob Schreiber Shail Aditya Bob Rau
Vinod Kathail Scott Mahlke
Darren Cronquist Mukund Sivaraman

HP Labs, Palo Alto

Outline

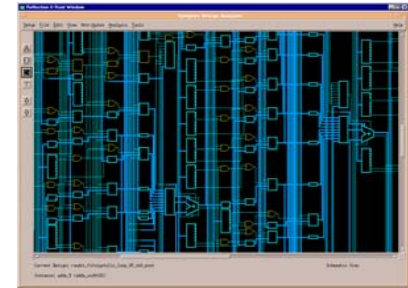
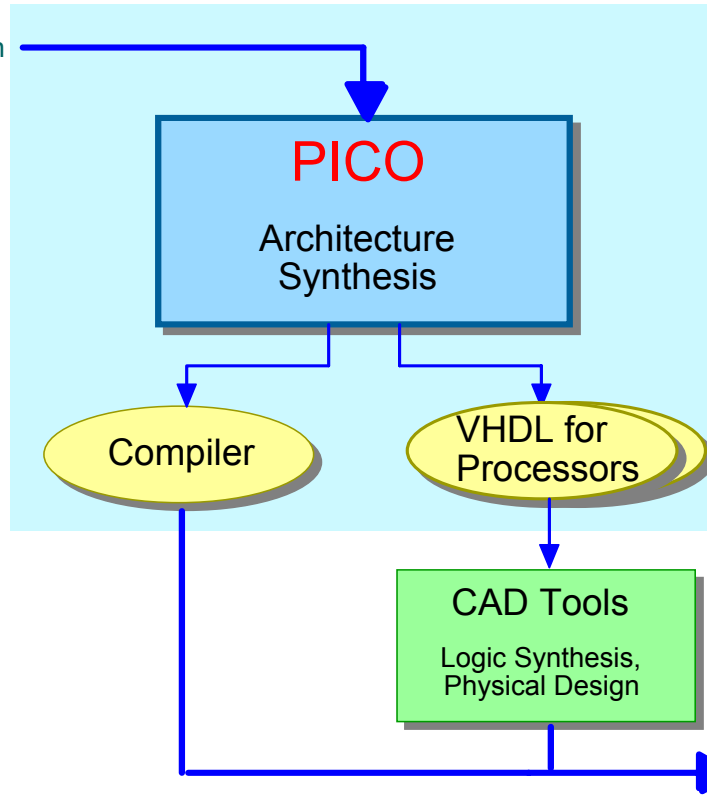
- What Can PICO Do for an SOC Designer?
- The PICO System Design Hierarchy
- From Sequential to Parallel Loop Nest
- Parallel Loop Nest to Processor Design

PICO overview

Program In

```
emacs@simpson.hpl.hp.com
Buffers Files Tools Edit Search C++ Help
#pragma bitsize outbuf 14

void
systolic_dct()
{
  int block, row_col;
  /* Do ROWS */
  for (block = 0; block < NBLOCKS; block++) {
    for (row_col = 0; row_col < 8; row_col++) {
      short in0, in1, in2, in3, in4, in5, in6, in7;
      -----Emacs: systolic_dct.c 9:34am (C++ CVS-1.1)
```

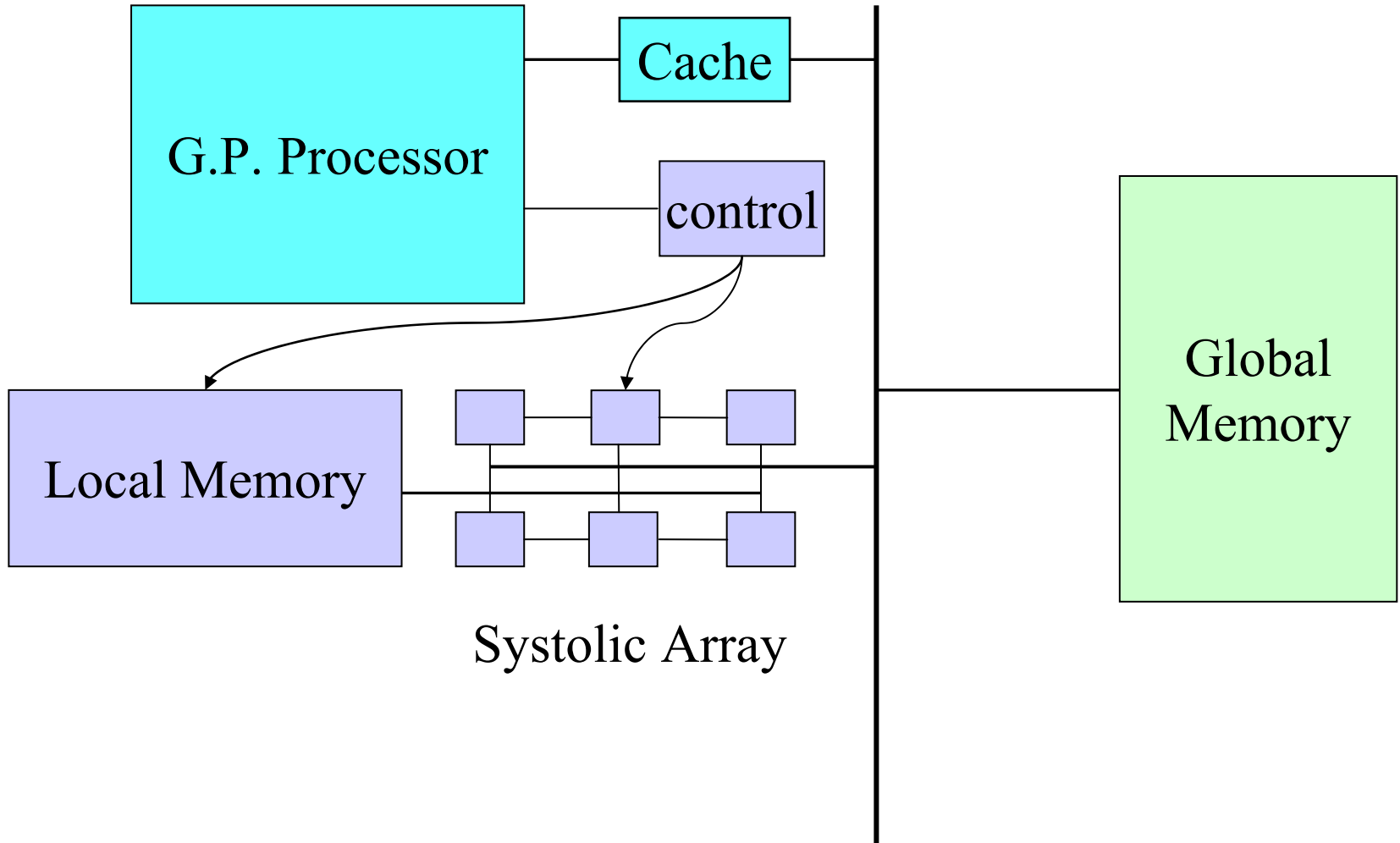


Program In --> IP Out

Using PICO

- User provides application, test data, and design space limits
- User indicates **hot** loop nests
- PICO creates Pareto set of ASIP designs.
- Each design has a customized VLIW with zero or more loop nests realized in HW
- User selects appropriate design for SOC based on area, power, performance tradeoff

PICO's ASIP Architecture



Hierarchical Design Frameworks

An Automated Design Template

Parameter
Ranges

Function
Specification

SpaceWalker

Constructor

Evaluator

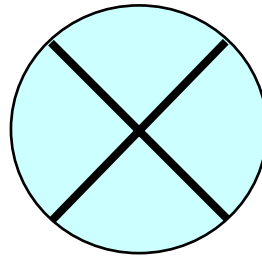
Pareto Filter

Good Systems from Good Subsystems

Cache
Pareto

VLIW
Pareto

NPA Pareto



System Constructor

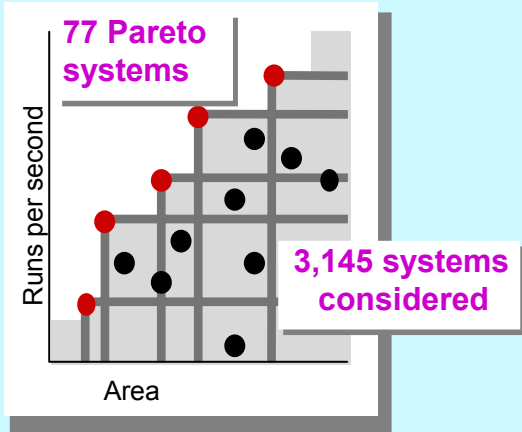
System Evaluator

System Pareto Filter

design space exploration

Design Space Exploration

Compile
Estimate Cycle Count

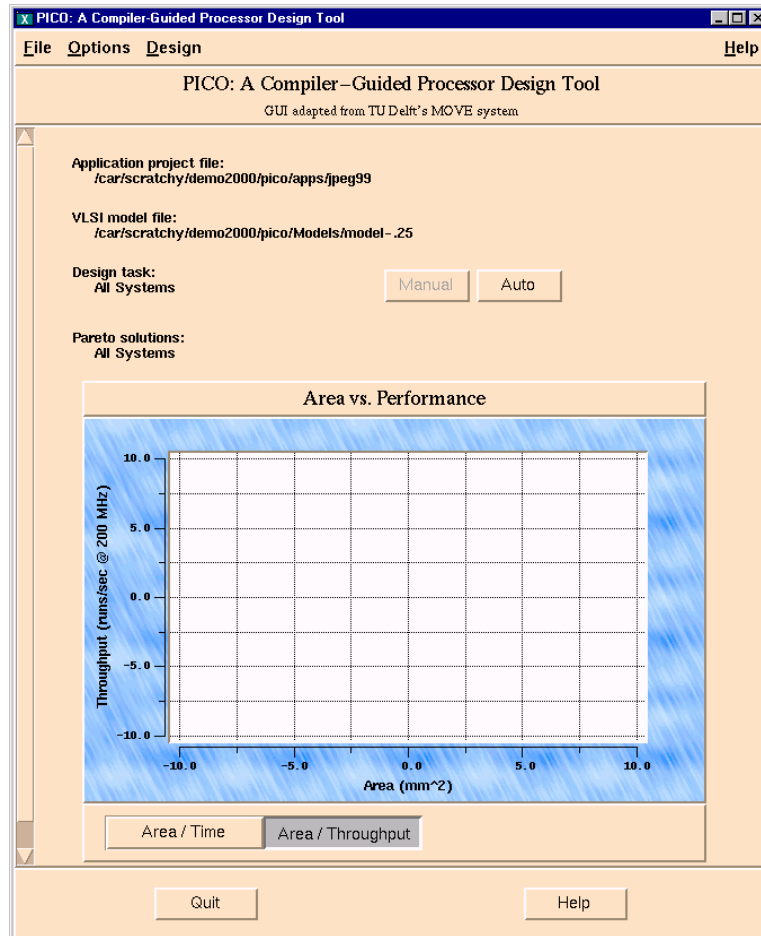


Synthesize
Estimate Area

2.5 million systems specified



PICO GUI



Limiting the Design Space

EPIC-Only + Hybrid Design Space Explorer: '/cat/scratchy/demo2000/.chipotle/data/jpeg99/transient/explore_range'

File Options Help

EPIC-Only + Hybrid Design Space Explorer

Predication: <input type="radio"/> general <input type="radio"/> none <input checked="" type="radio"/> both	Speculation: <input checked="" type="radio"/> general <input type="radio"/> restricted <input type="radio"/> both	EPIC Type: <input checked="" type="radio"/> IFMB <input type="radio"/> heterogeneous	Systolic: <input type="radio"/> with <input type="radio"/> without <input checked="" type="radio"/> both
---	---	---	--

Functional Units

Integer:	low: <input type="text" value="1"/>	high: <input type="text" value="8"/>	step: <input type="text" value="1"/>
Float:	low: <input type="text" value="1"/>	high: <input type="text" value="1"/>	step: <input type="text" value="1"/>
Memory:	low: <input type="text" value="1"/>	high: <input type="text" value="4"/>	step: <input type="text" value="1"/>
Branch:	low: <input type="text" value="1"/>	high: <input type="text" value="1"/>	step: <input type="text" value="1"/>

Register Files

Integer:	low: <input type="text" value="32"/>	high: <input type="text" value="128"/>	step: <input type="text" value="32"/>
Float:	low: <input type="text" value="32"/>	high: <input type="text" value="32"/>	step: <input type="text" value="16"/>
Predicate:	low: <input type="text" value="32"/>	high: <input type="text" value="32"/>	step: <input type="text" value="16"/>
Branch:	low: <input type="text" value="16"/>	high: <input type="text" value="16"/>	step: <input type="text" value="16"/>

Systolic Array

Systolic Processors:	low: <input type="text" value="1"/>	high: <input type="text" value="2"/>	step: <input type="text" value="1"/>
II:	low: <input type="text" value="1"/>	high: <input type="text" value="8"/>	step: <input type="text" value="1"/>
Memory Ports:	low: <input type="text" value="1"/>	high: <input type="text" value="1"/>	step: <input type="text" value="1"/>

Level 1 data cache parameters

Cache sets:	low: <input type="text" value="128"/>	high: <input type="text" value="128"/>	step: <input type="text" value="1"/>
Associativity:	low: <input type="text" value="2"/>	high: <input type="text" value="2"/>	step: <input type="text" value="1"/>
Line size:	low: <input type="text" value="16"/>	high: <input type="text" value="32"/>	step: <input type="text" value="1"/>

Level 1 instruction cache parameters

Cache sets:	low: <input type="text" value="128"/>	high: <input type="text" value="128"/>	step: <input type="text" value="1"/>
Associativity:	low: <input type="text" value="2"/>	high: <input type="text" value="2"/>	step: <input type="text" value="1"/>
Line size:	low: <input type="text" value="32"/>	high: <input type="text" value="128"/>	step: <input type="text" value="1"/>

Level 2 unified cache parameters

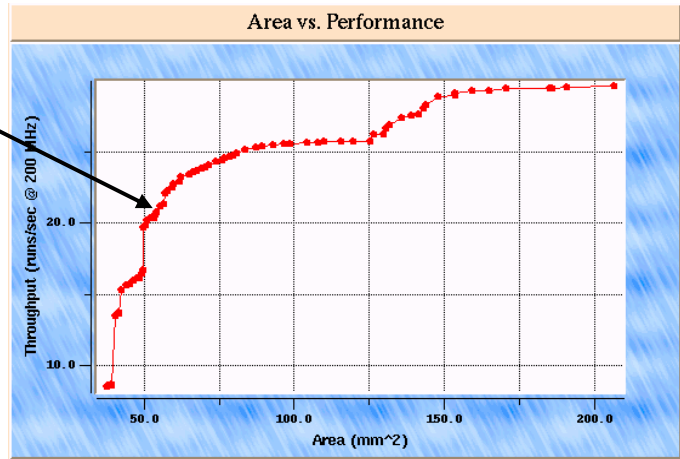
Cache sets:	low: <input type="text" value="256"/>	high: <input type="text" value="256"/>	step: <input type="text" value="1"/>
Associativity:	low: <input type="text" value="2"/>	high: <input type="text" value="3"/>	step: <input type="text" value="1"/>
Line size:	low: <input type="text" value="64"/>	high: <input type="text" value="128"/>	step: <input type="text" value="1"/>

Search mode: Show final Pareto Show intermediate Paretos

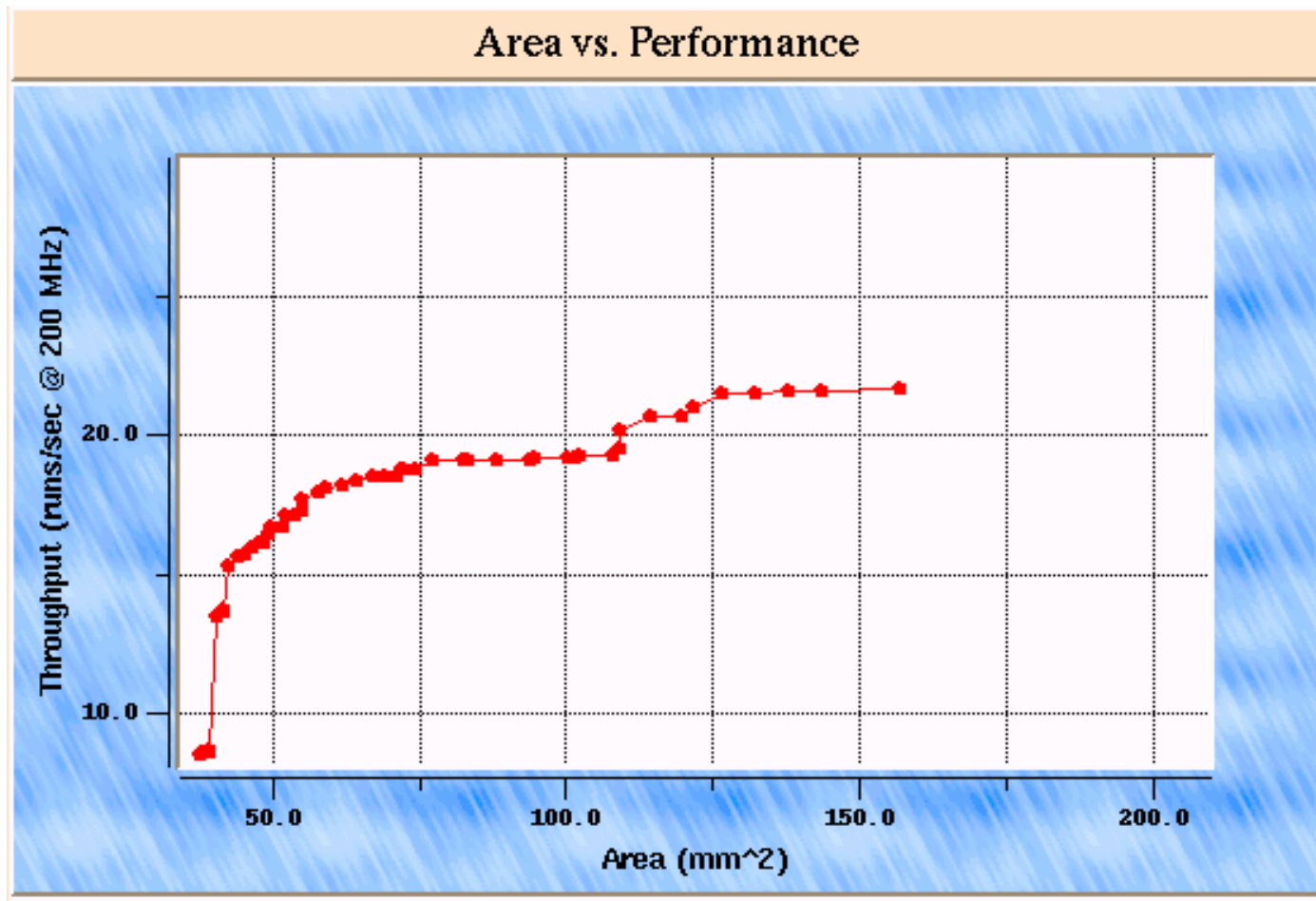
Explore Quit Help

Exploration

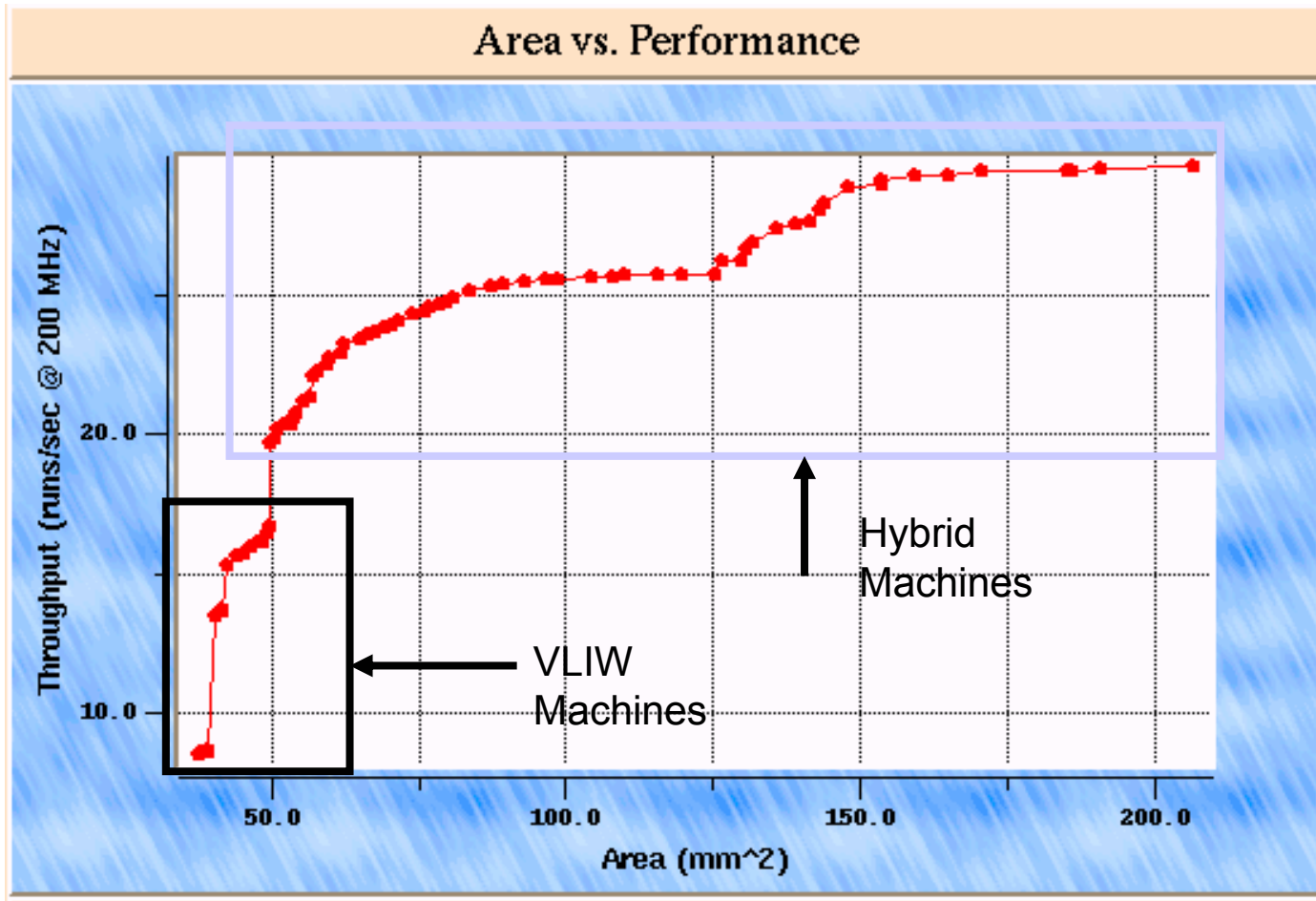
```
hpterm [simpson via TELNET]
>> [systemHandle = 15861, expHandle = 15871]
>> area = 68.5671, cycles = 9.81964e+06, packetSize = 512, textSize = 105624 (0.828125)
>> Evaluated jpeg99 on vliw(p0s1 i6f1m4b1 i[96,0]f[32,0]p[32,0]b[16] IFMT):
>> [systemHandle = 15881, expHandle = 15891]
>> area = 72.2403, cycles = 9.8348e+06, packetSize = 512, textSize = 98728 (0.828125)
>> Added vliw(p1s1 i7f1m4b1 i[64,0]f[32,0]p[32,0]b[16] IFMT) 15901
>> Added vliw(p0s1 i8f1m4b1 i[64,0]f[32,0]p[32,0]b[16] IFMT) 15921
>> Added vliw(p0s1 i7f1m4b1 i[96,0]f[32,0]p[32,0]b[16] IFMT) 15941
>> STEP 11: 66 machines evaluated
>> Selected vliw(p1s1 i7f1m4b1 i[64,0]f[32,0]p[32,0]b[16] IFMT)
>> Selected vliw(p0s1 i8f1m4b1 i[64,0]f[32,0]p[32,0]b[16] IFMT)
>> Selected vliw(p0s1 i7f1m4b1 i[96,0]f[32,0]p[32,0]b[16] IFMT)
>> Evaluated jpeg99 on vliw(p1s1 i7f1m4b1 i[64,0]f[32,0]p[32,0]b[16] IFMT):
>> [systemHandle = 15901, expHandle = 15911]
>> area = 70.946, cycles = 1.05966e+07, packetSize = 512, textSize = 132284 (1)
>> Evaluated jpeg99 on vliw(p0s1 i8f1m4b1 i[64,0]f[32,0]p[32,0]b[16] IFMT):
>> [systemHandle = 15921, expHandle = 15931]
>> area = 75.6306, cycles = 8.92249e+06, packetSize = 512, textSize = 115676 (1)
>> Evaluated jpeg99 on vliw(p0s1 i7f1m4b1 i[96,0]f[32,0]p[32,0]b[16] IFMT):
>> [systemHandle = 15941, expHandle = 15951]
>> area = 79.9293, cycles = 9.81246e+06, packetSize = 512, textSize = 106200 (0.828125)
>> Added vliw(p1s1 i8f1m4b1 i[64,0]f[32,0]p[32,0]b[16] IFMT) 15961
>> Added vliw(p0s1 i8f1m4b1 i[96,0]f[32,0]p[32,0]b[16] IFMT) 15981
>> STEP 12: 69 machines evaluated
>> Selected vliw(p1s1 i8f1m4b1 i[64,0]f[32,0]p[32,0]b[16] IFMT)
>> Selected vliw(p0s1 i8f1m4b1 i[96,0]f[32,0]p[32,0]b[16] IFMT)
>> Evaluated jpeg99 on vliw(p1s1 i8f1m4b1 i[64,0]f[32,0]p[32,0]b[16] IFMT):
>> [systemHandle = 15961, expHandle = 15971]
>> area = 77.2607, cycles = 1.05966e+07, packetSize = 512, textSize = 141764 (1.19531)
>> Evaluated jpeg99 on vliw(p0s1 i8f1m4b1 i[96,0]f[32,0]p[32,0]b[16] IFMT):
>> [systemHandle = 15981, expHandle = 15991]
>> area = 88.8775, cycles = 8.91315e+06, packetSize = 512, textSize = 112832 (1)
>> Added vliw(p1s1 i8f1m4b1 i[96,0]f[32,0]p[32,0]b[16] IFMT) 16001
>> Added vliw(p0s1 i8f1m4b1 i[128,0]f[32,0]p[32,0]b[16] IFMT) 16021
>> STEP 13: 71 machines evaluated
>> Selected vliw(p1s1 i8f1m4b1 i[96,0]f[32,0]p[32,0]b[16] IFMT)
>> Selected vliw(p0s1 i8f1m4b1 i[128,0]f[32,0]p[32,0]b[16] IFMT)
>> Evaluated jpeg99 on vliw(p1s1 i8f1m4b1 i[96,0]f[32,0]p[32,0]b[16] IFMT):
>> [systemHandle = 16001, expHandle = 16011]
>> area = 90.7583, cycles = 1.059e+07, packetSize = 512, textSize = 144724 (1.19531)
>> Evaluated jpeg99 on vliw(p0s1 i8f1m4b1 i[128,0]f[32,0]p[32,0]b[16] IFMT):
>> [systemHandle = 16021, expHandle = 16031]
>> area = 100.992, cycles = 8.92426e+06, packetSize = 512, textSize = 111312 (1)
>> STEP 14: 73 machines evaluated
>> STEP 15: 73 machines evaluated
>> Composing (vliw + memory) paretos.....
>> done.
>> Refining system cost estimates...
>> done.
```



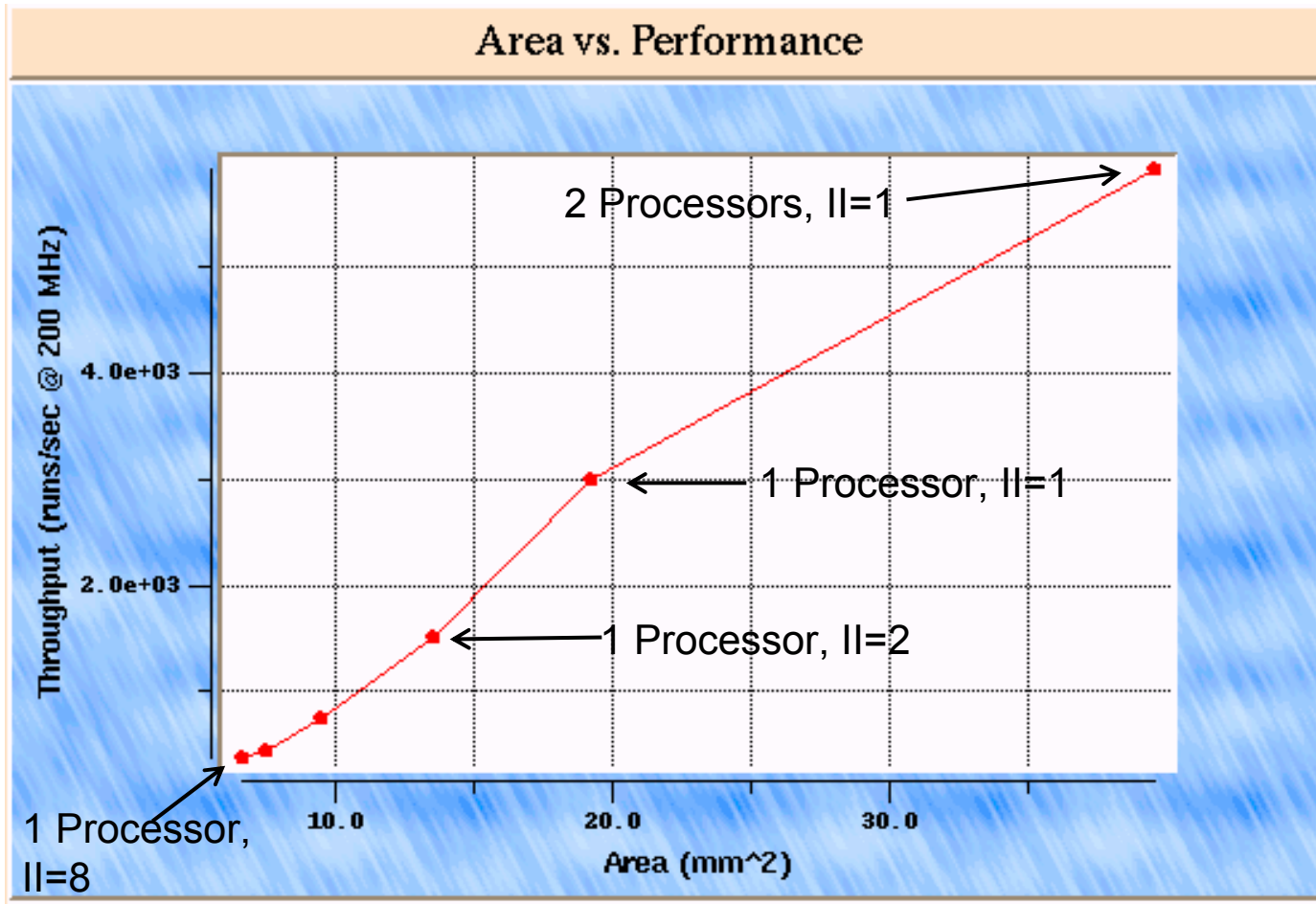
Pareto Optimal Machines: VLIW-only



Pareto Optimal Machines: All systems



Systemic Design: Exploration



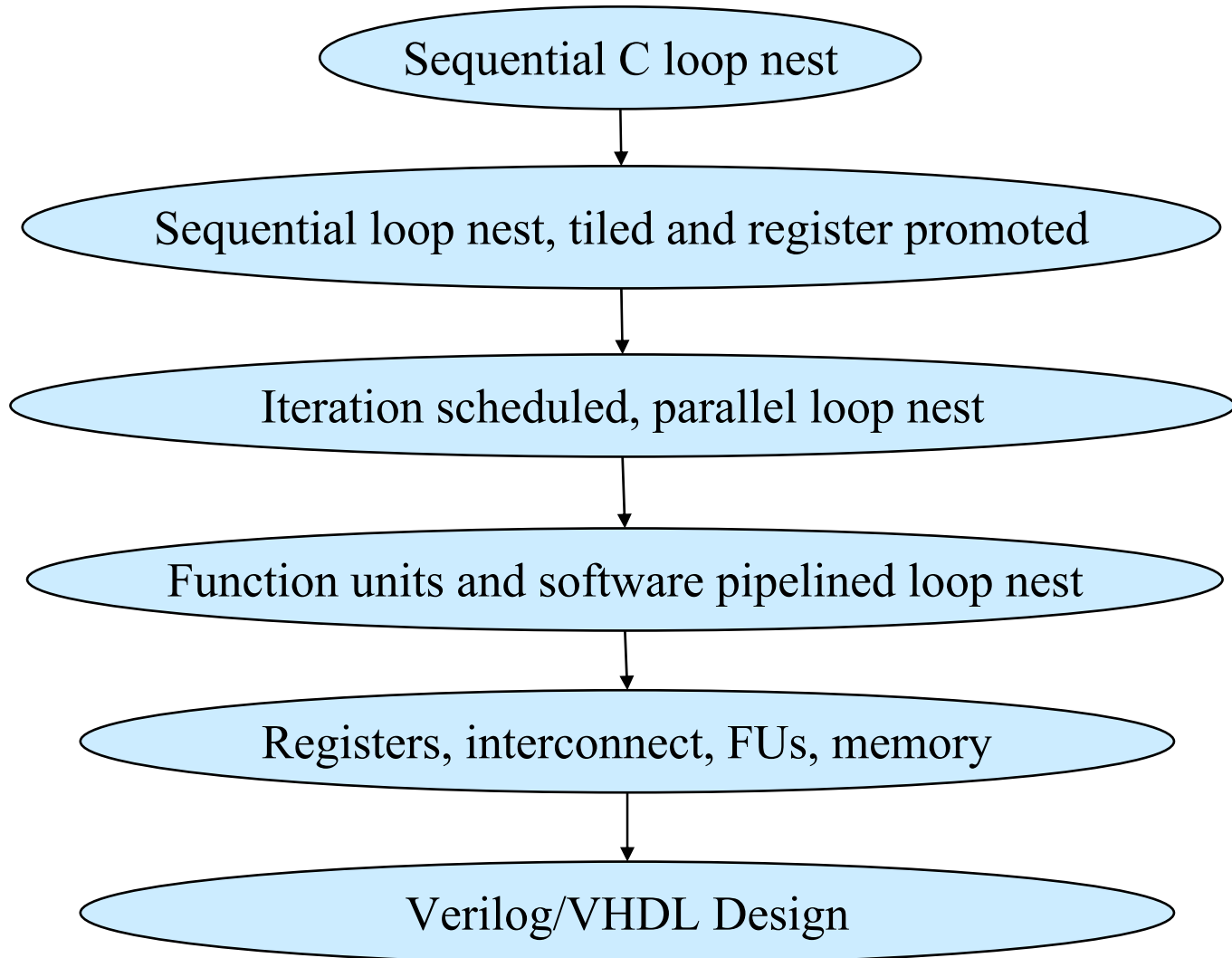
Synthesis of a Non-Programmable, Application-Specific Accelerator:

From Sequential Loop Nest to Parallel Loop Nest

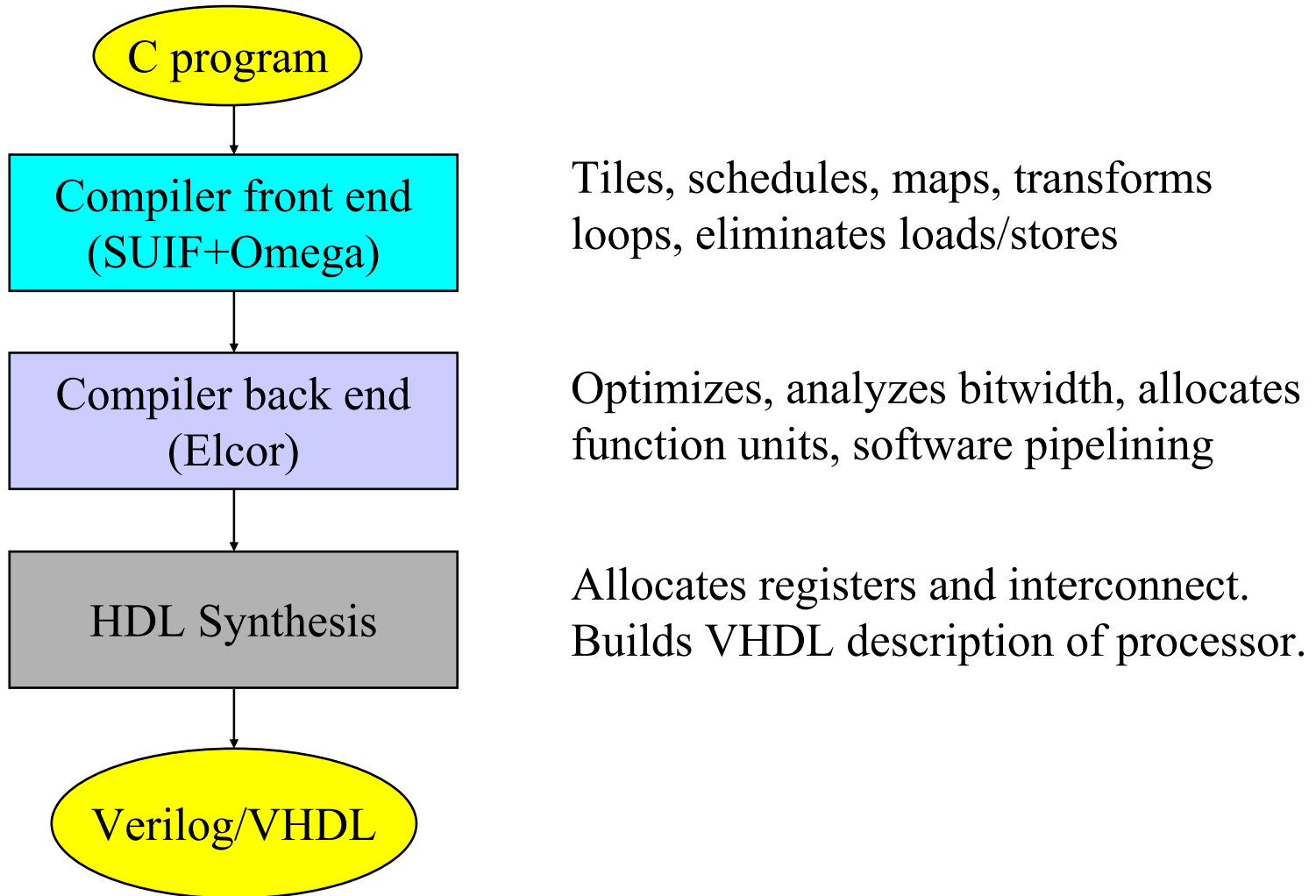
Input Language

- A perfect loop nest \rightarrow A systolic array
- A sequence of nests \rightarrow A pipeline of arrays
- Constant loop bounds
- Dependence analysis must be feasible:
 - No aliasing through pointers
- Language extensions
 - `#pragma bitsize x 12`
 - `#internal coeff`

From C to VHDL

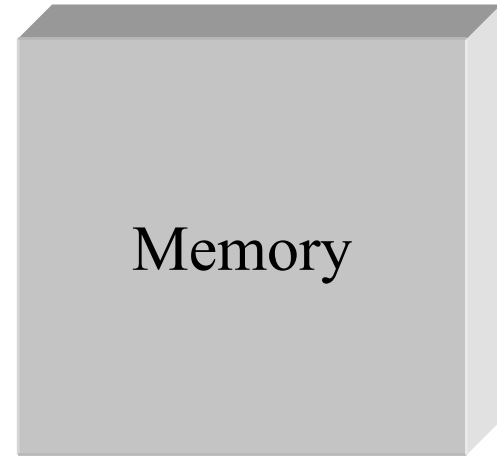
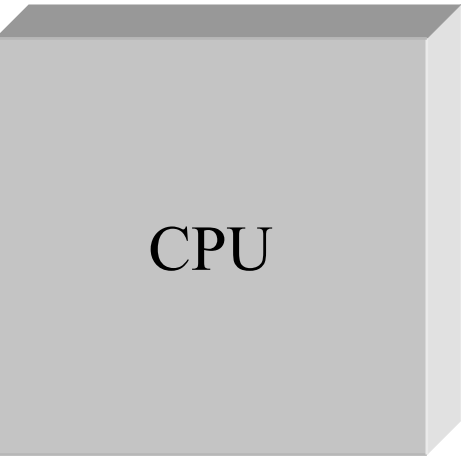


From C to VHDL

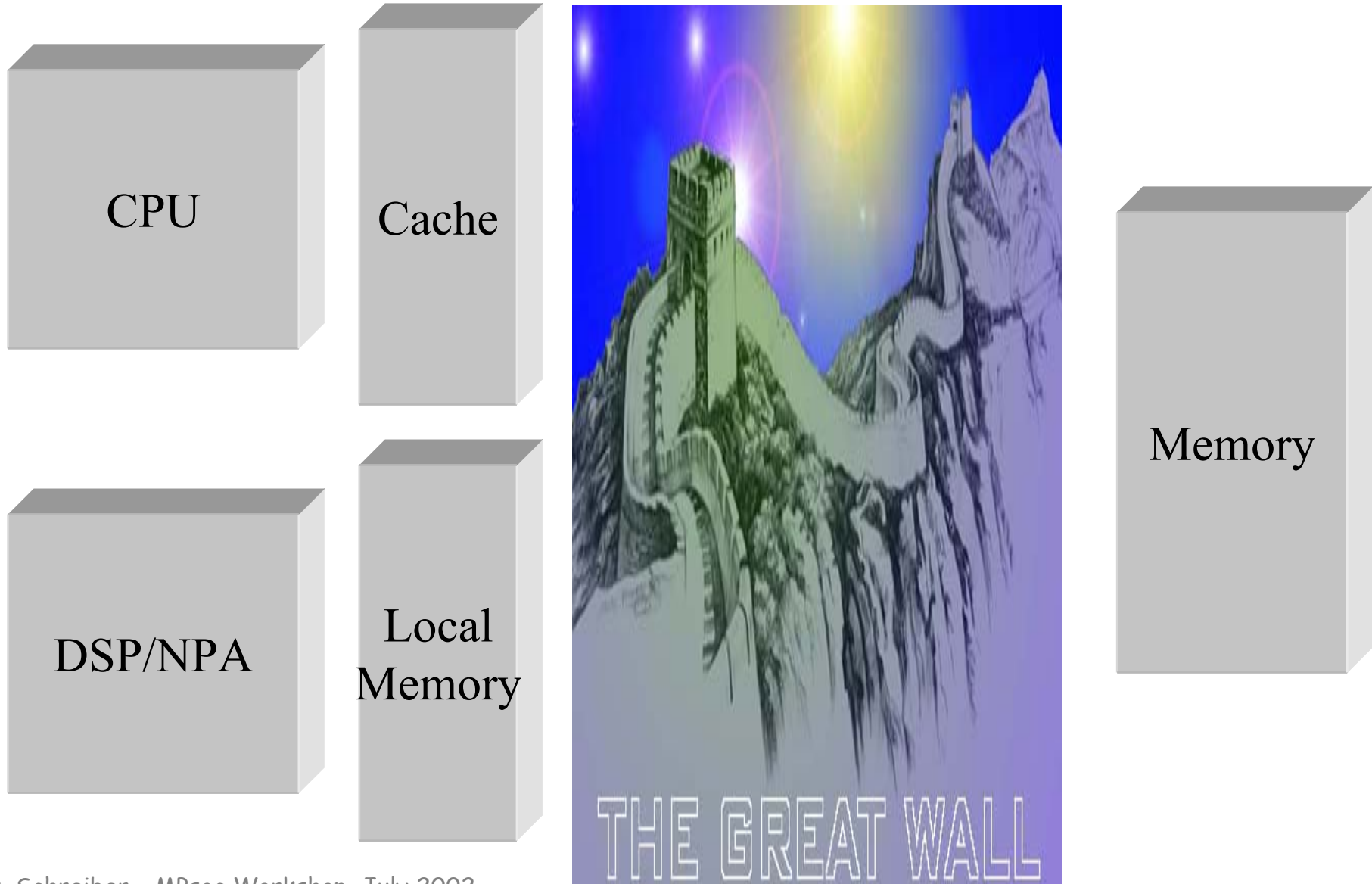


What does it take to make this
efficient?

The Memory Wall



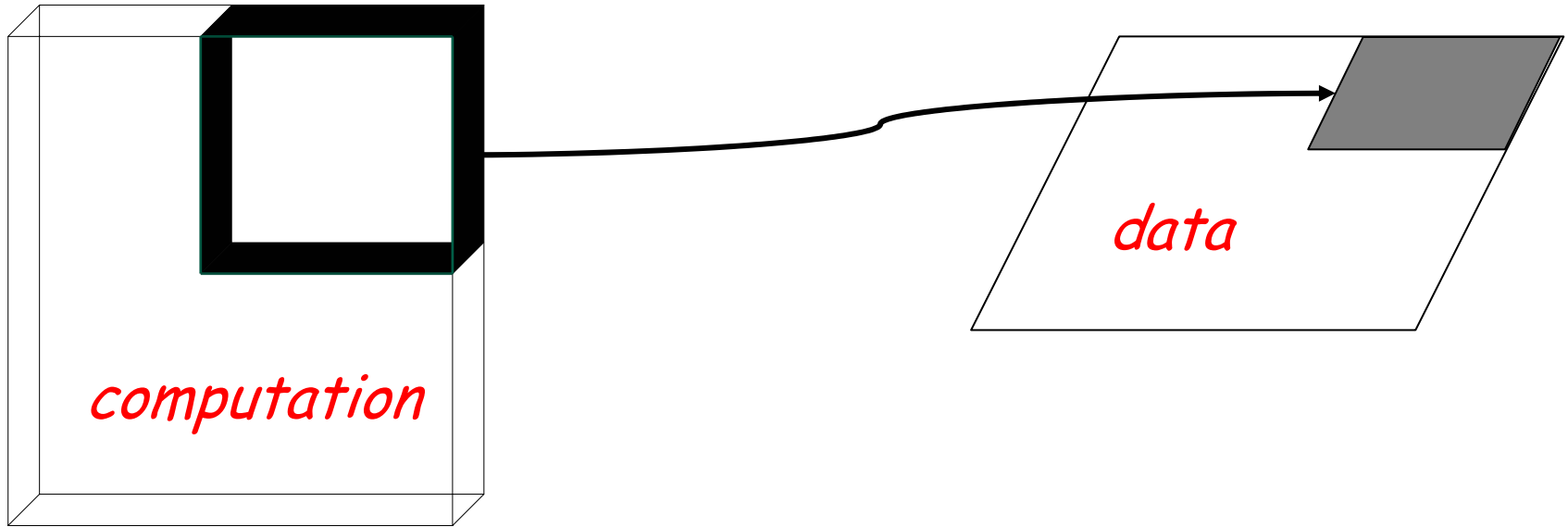
Cache and Local Memory



Goal of Code Transformation

```
for each TILE
{
  for (t = 0; t < Tfinal; t++)
  {
    forall processors p
    {
      X[t][p] = . . .
      Y[t-1][p+1] . . .
    }
  }
}
```

Tiling the Iteration Space



Volume/Surface = $O(\text{radius})$

Computation/Footprint = $\Omega(\text{radius})$

Computation/Footprint = CPU/Memory

Load/Store Elimination

- For affine array references, intermediate results in registers
- For affine, read-only array references, data routed through registers; no value loaded more than once.

Tile Shapes

Big tiles → More local memory

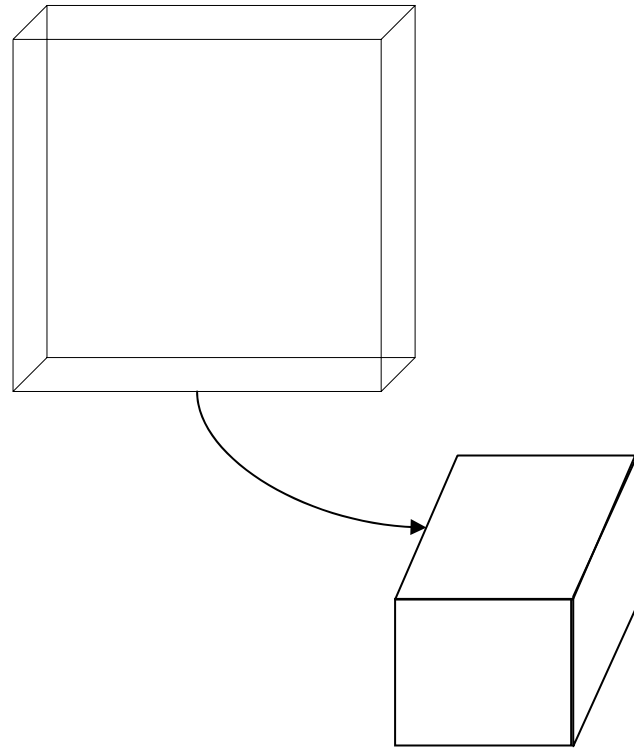
Small tiles → less reuse of data, more global memory bandwidth

Optimal tile → smallest tile that does not oversubscribe memory bandwidth

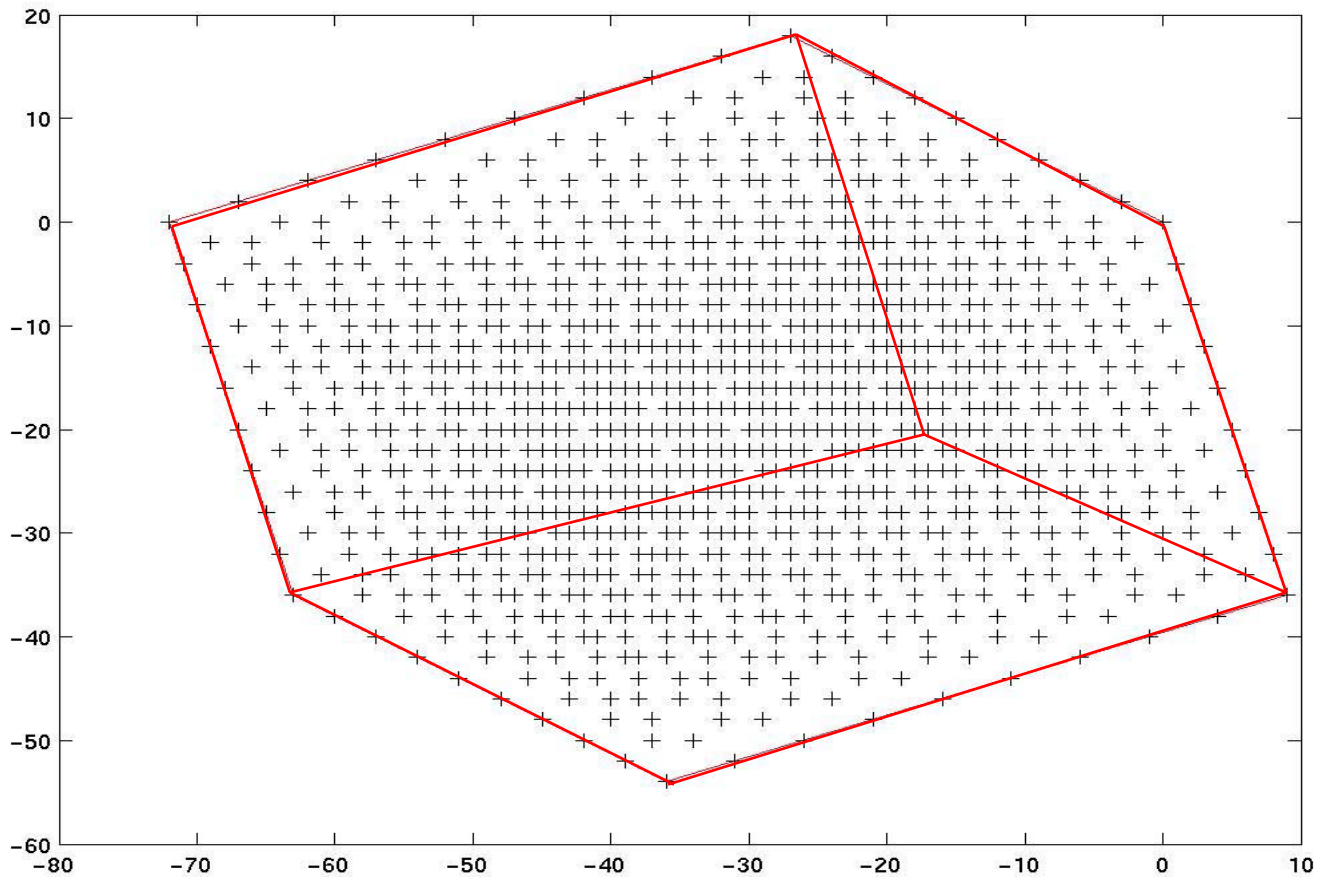
Estimating the Footprint

Affine array reference
 $X[i+j][2*j-3*k]$

How many integer points in
an affine image of a
rectangular iteration space?



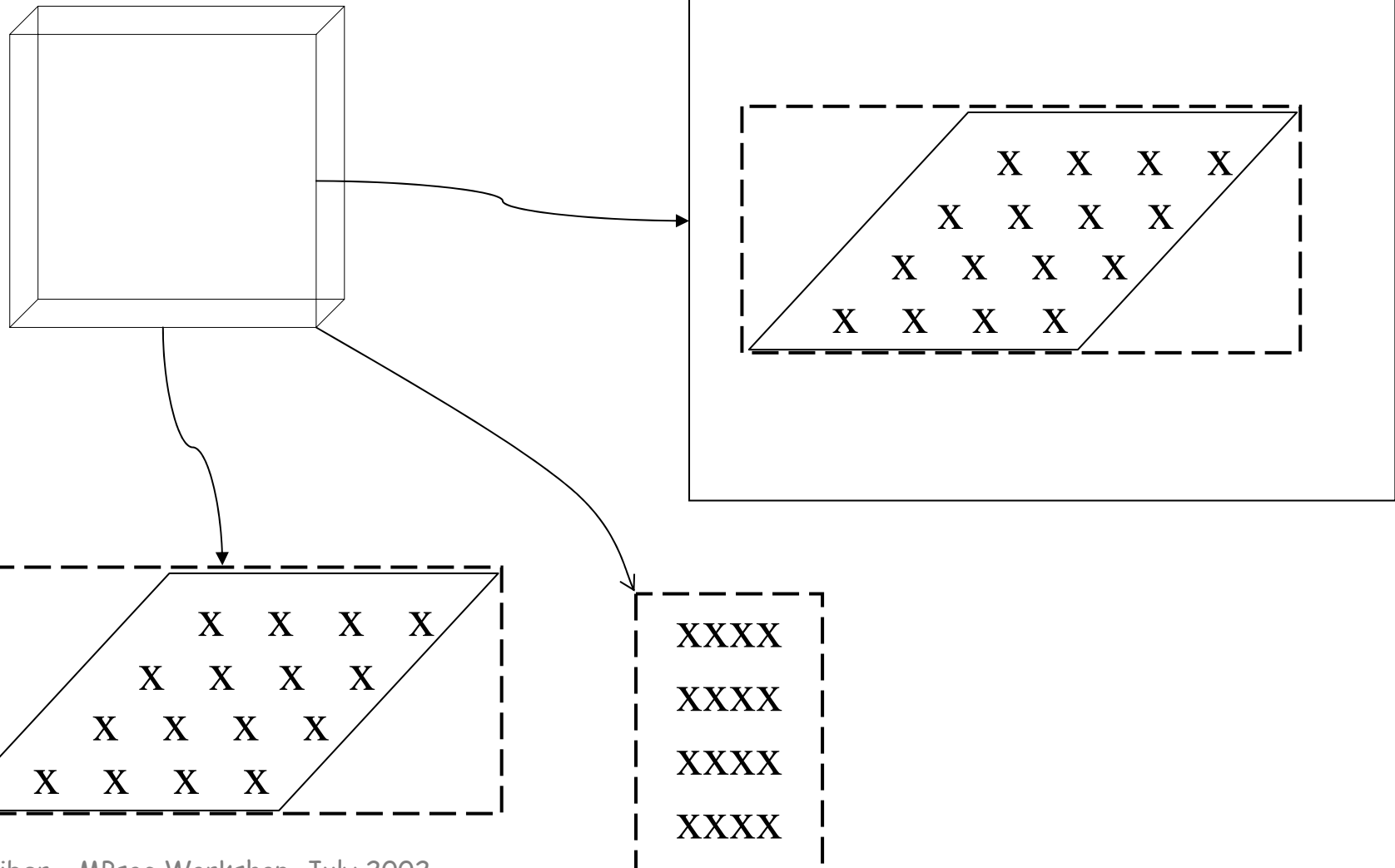
Example: the Affine Image of an Iteration Space



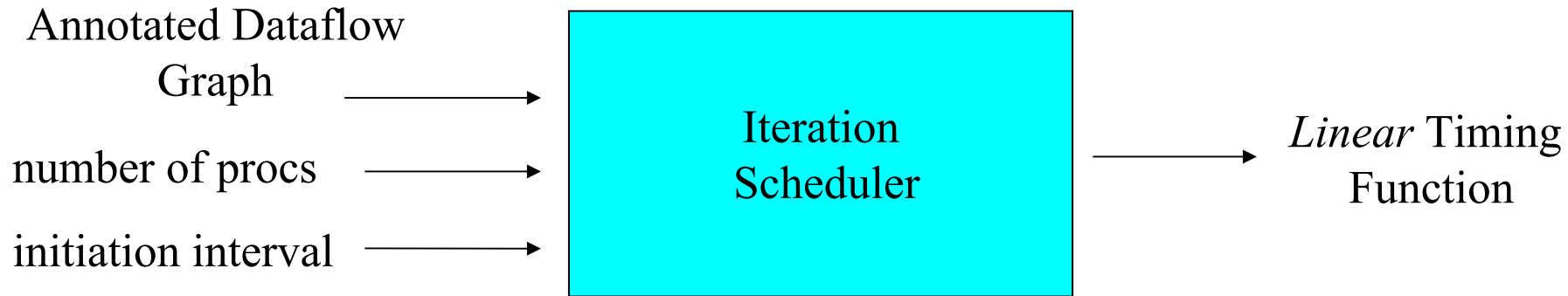
Corrected Estimates

- Published bounds on the size of the image of a Z -polytope are wrong
- Our corrections:
 - footprint = iteration space for 1-1 mappings
 - 1-1 if no integer null vector in the iteration space
 - corrected bounds from finding number of iterations that differ by a null vector
 - within 20 percent in practice

Reindexing to Reduce Local Memory



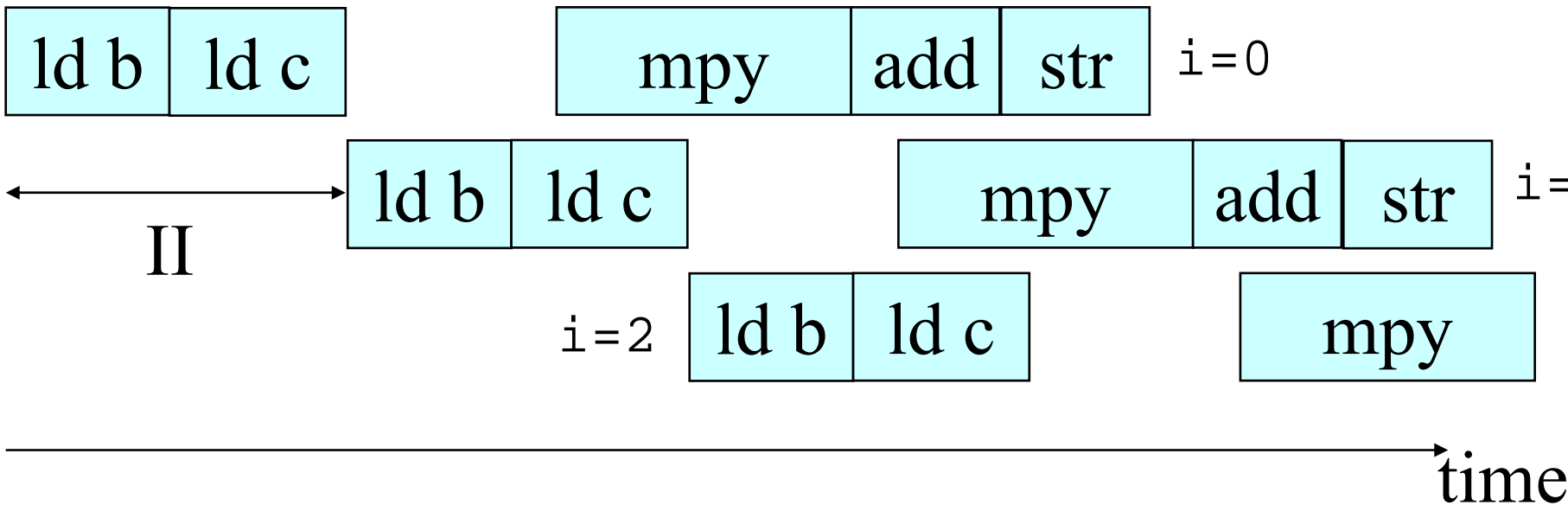
Finding the Parallel Iteration Schedule



- *Processors* a mesh of processors is given
- *Initiation Interval (II)* every processor starts an iteration periodically with period equal to II (*hardware* pipelining)
- *Mapping* clusters of iterations are mapped to each processor
- *Schedule* one iteration per processor every II cycles
- *Honor* data dependence constraints
- *Find the schedule* via efficient direct search method

Hardware/Software Pipelining

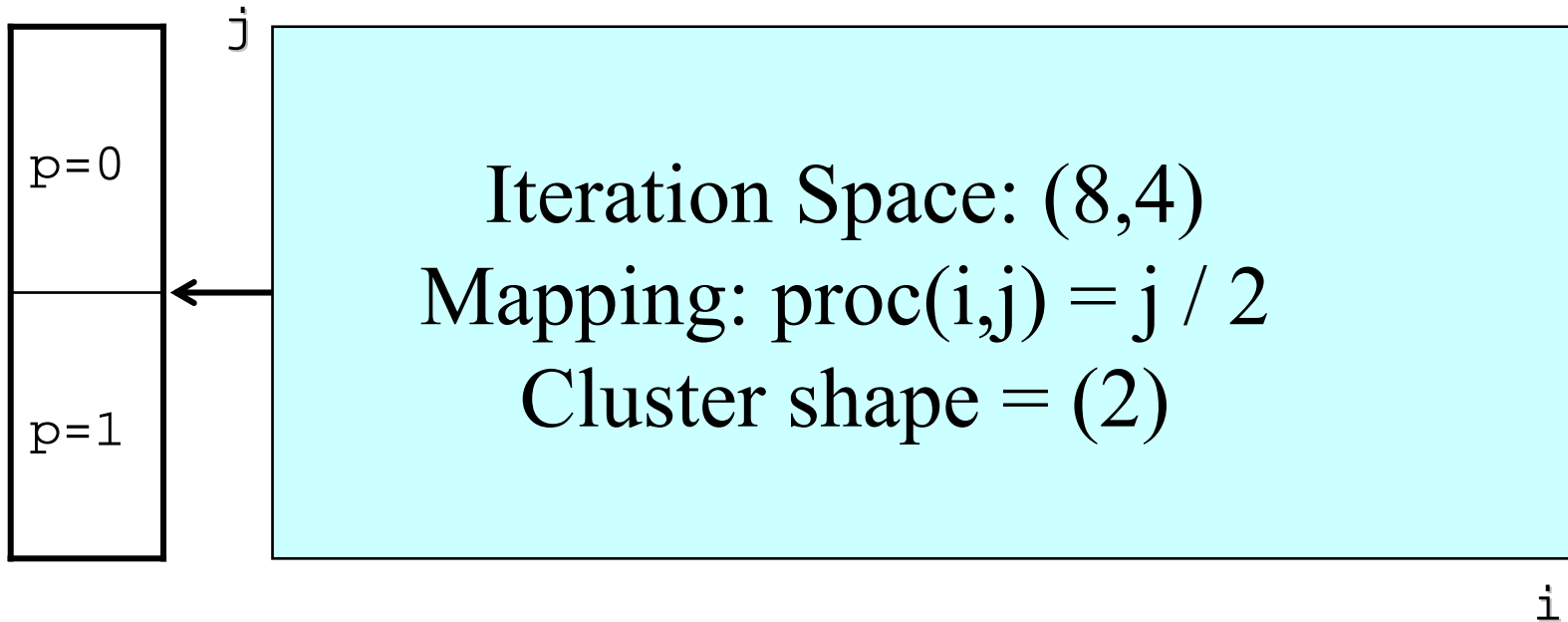
```
for (i=0; i < 100; i++) a[i] += b[i]*c[i]
```



Lower Bounds on II (RecMII, ResMII)

The Mapping of Iterations to Processors

```
for (i = 0; i < 8; i++)  
  for (j = 0; j < 4; j++)  
  {  
    y[i] += w[j] * x[i-j];  
  }
```



A Tight Schedule: $(i,j) \rightarrow 2i+3j$

```
for (i = 0; i < 8; i++)  
  for (j = 0; j < 4; j++)  
  {  
    y[i] += w[j] * x[i-j];  
  }
```

	j	9	11	13	15	17	19	21	23
p=0		6	8	10	12	14	16	18	20
p=1		3	5	7	9	11	13	15	17
		0	2	4	6	8	10	12	14

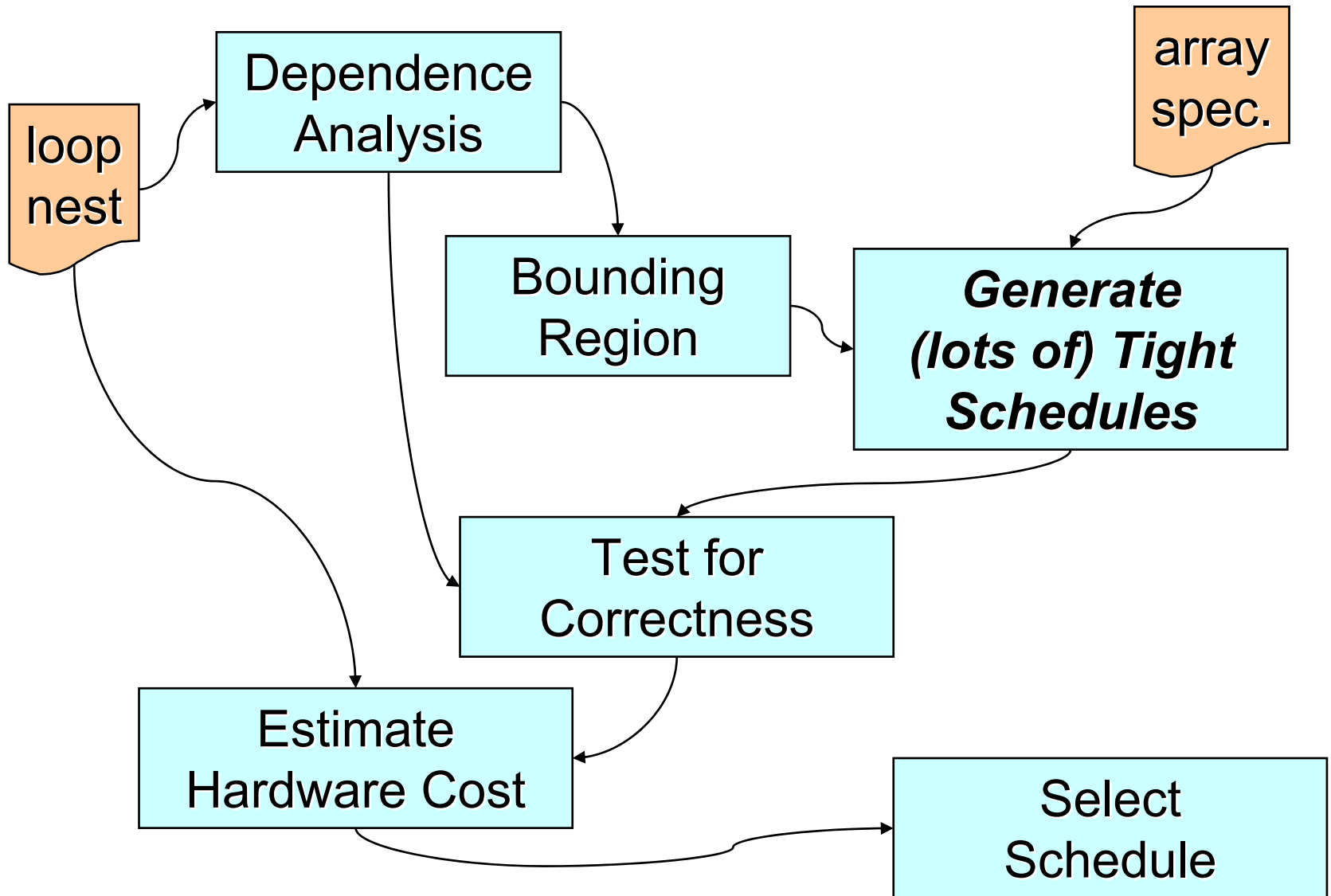
Diagram illustrating a tight schedule for the nested loops. The vertical axis is labeled j and the horizontal axis is labeled i . The schedule is divided into two parts, $p=0$ and $p=1$. The values in the table represent the expression $2i+3j$. Blue arrows indicate dependencies: a vertical arrow from 0 to 3, a diagonal arrow from 0 to 5, and a horizontal arrow from 0 to 2. A black arrow points from the $p=1$ section to the $p=0$ section.

Tight Schedules – Prior Work

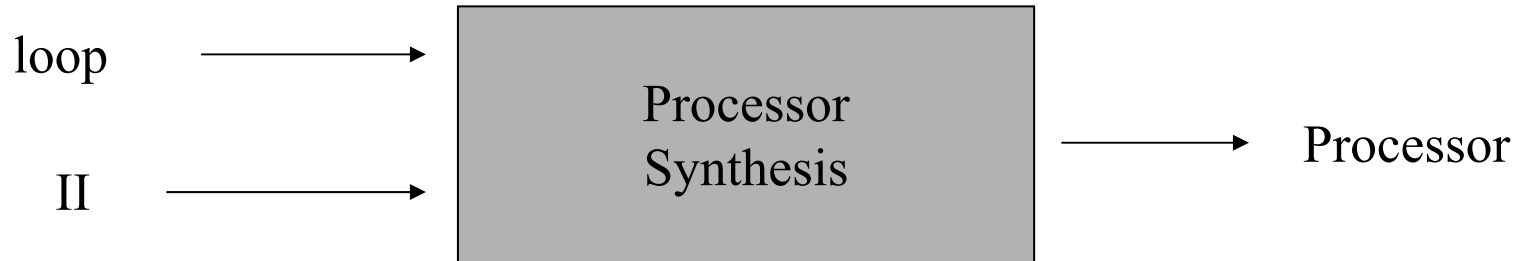
Darte/Delosme, Chen/Megson.

- *GIVEN*: Iteration space, projection direction, linear schedule
- *DETERMINE*: The allowed cluster shapes
- *Tail Wags Dog!*

Constructing the Schedule



Processor Synthesis

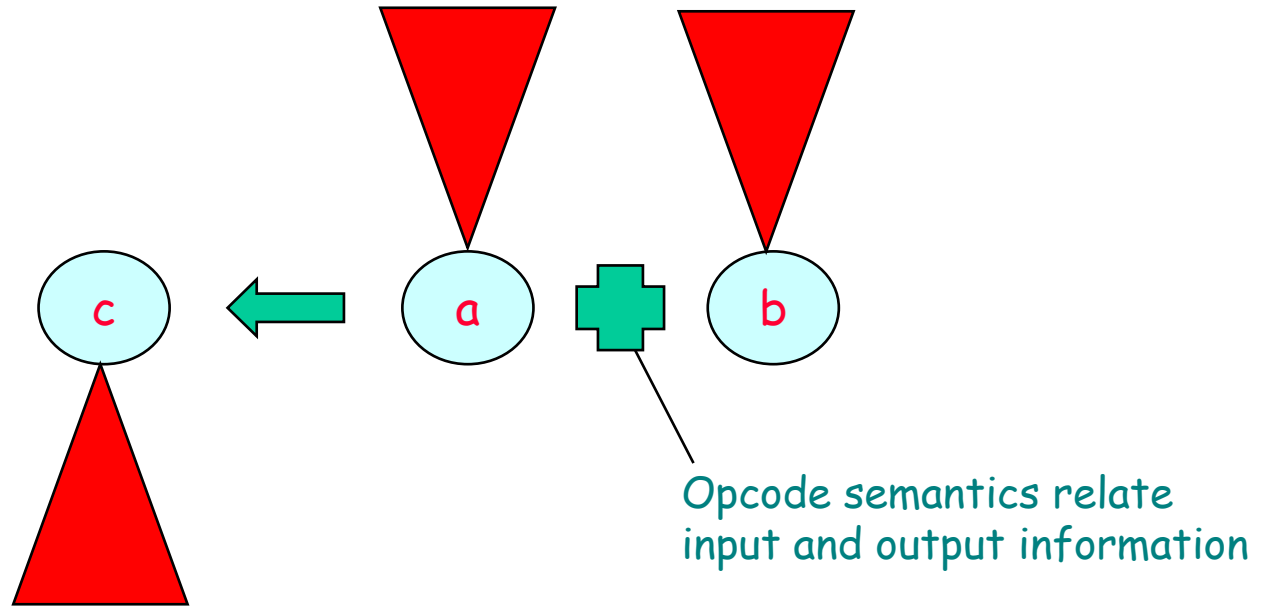


- *Optimize* the loop body
- *Analyze* bitwidth of all values
- *Allocate* the function units
- *Map* operations to function units
- *Schedule* operations
- *Allocate* registers and memory
- *Interconnect* communicating elements

Parallel, custom, designed to spec: EFFICIENT!

Bitwidth analysis - basic idea

Input information limits the amount of information that can be produced



Information required by consumers limits the amount that must be produced

Optimal FU allocation

count	FU cost	type	Operation type	count
1	10	+	+	3
0	10	-	-	1
1	13	+/-		

The diagram shows three FU types on the left and two operation types on the right. Lines connect the FU types to the operation types with weights:

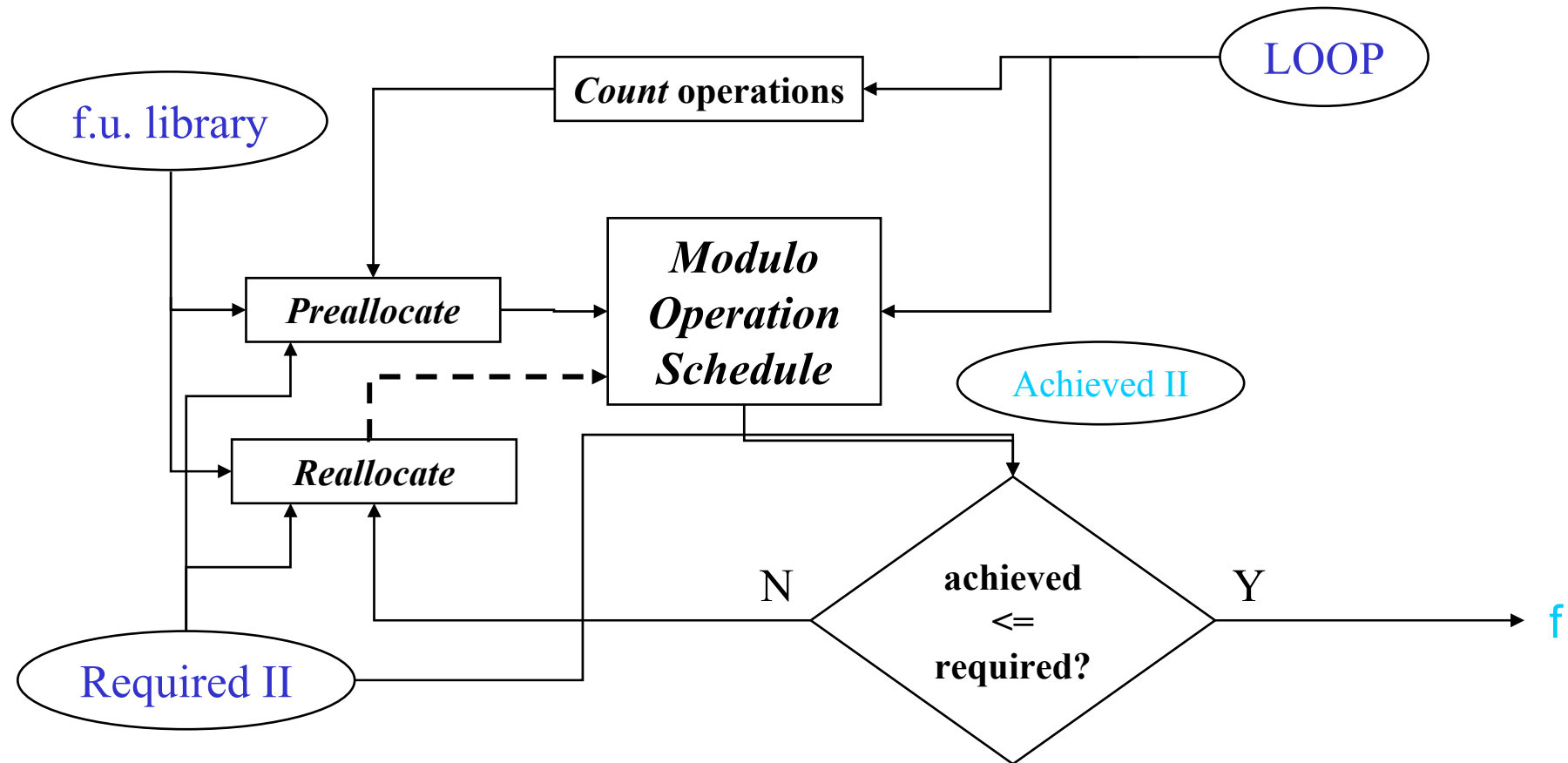
- A line connects FU type '+' (cost 10) to operation type '+' (count 3) with a weight of 2.
- A line connects FU type '-' (cost 10) to operation type '-' (count 1) with a weight of 1.
- A line connects FU type '+/-' (cost 13) to operation type '-' (count 1) with a weight of 1.

MILP: minimize cost subject to sufficient capacity

Allocation and Op Scheduling

Given: Inner loop and II

Find: Cheapest processor that achieves II on the loop



Conclusions

- Accurate static analysis of memory bandwidth – optimal tiling
- Linear iteration scheduling: solved problem
- Efficient datapath synthesis – a hard problem, good heuristics
- Automatic NPA synthesis is practical
- Automatic synthesis of full embedded systems is feasible, too

Related publications :

Robert Schreiber, Shail Aditya, Scott Mahlke, Vinod Kathail, B. Ramakrishna Rau, Darren Cronquist, and Mukund Sivaraman.

PICO-NPA: High-level synthesis of nonprogrammable hardware accelerators.
In Journal of VLSI Signal Processing 31: 127-142 (2002).

Shail Aditya, B. Ramakrishna Rau, and Vinod Kathail.

Automatic architecture synthesis of VLIW and EPIC processors.

In Proceedings of the 12th International Symposium on System Synthesis, San Jose, California, pp. 107--113, November 1999.

Alain Darte, Robert Schreiber, B. Ramakrishna Rau, and Frederic Vivien.

Constructing and exploiting linear schedules with prescribed parallelism.

ACM Transactions on Design Automation for Electronic Systems, 7(1), (2002)

Kyle Gallivan, William Jalby, and Dennis Gannon.

On the problem of optimizing data transfers for complex memory systems.

In Proceedings of the 1988 ACM International Conference on Supercomputing, pp. 238--253, 1988.

Scott Mahlke, Rajiv Ravindran, Michael Schlansker, Robert Schreiber, and Timothy Sherwood.

Bitwidth cognizant architecture synthesis of custom hardware accelerators.

IEEE Transactions on Computer-Aided Design of Circuits and Systems, 20(10):1-17, 2001.

William Pugh.

The Omega test: a fast and practical integer programming algorithm for dependence analysis.

Communications of the ACM, 35(8):102--114, 1992.

Patrice Quinton and Yves Robert.

Systolic Algorithms and Architectures.

Prentice Hall International (UK) Ltd., Hemel Hempstead, England, 1991.

B. Ramakrishna Rau.

Iterative modulo scheduling.

International Journal of Parallel Processing, 24:3--64, 1996.

B. Ramakrishna Rau, Vinod Kathail, and Shail Aditya.

Machine-description driven compilers for EPIC and VLIW processors.

Design Automation for Embedded Systems, 4:71--118, 1999.