

```

short a[100], b[100];
short dot-product (void )
{
  short dot = 0;
  for (int i = 0; i < 100; i++) {
    dot = dot + a[i] * b[i];
  }
  return dot;
}

```

Dot product (C code)

```

movl a,R0 || movl b,R1
clr A0 || bkrp 99,LBL
ld (R0)+,D2 || ld (R1)+,D3
nop
mac D2,D3,A0
nop
LBL: ret

```

Dot product (assembly code)

To solve this problem several approaches at different stages of the compilation flow are possible :

- ▶ High level intermediate representation
- ▶ Low level intermediate representation

Variable Partitioning

Viera Šipkova
 Christian Doppler Laboratory
 Compilation Techniques for Embedded Processors
 Institute for Computer Languages
 Vienna University of Technology

Efficient Variable Allocation to dual Memory Banks of DSPs

- ▶ Motivation
- ▶ Introduction
- ▶ Partitioning Scheme
- ▶ Experimental Results
- ▶ Conclusion

Outline

To improve the overall performance -

- ▶ DSPs are equipped with *multiple data memory banks*.
- ▶ The compiler must *partition program variables* into memory banks.

Motivation

Interference Graph

The edge-weighted undirected graph $G = (V, E)$, where each vertex $v \in V$ – represents a memory access – indicates the possible parallel access to v and w (different memory banks)

edge $e = (v, w) \in E$ –

Interference Edge

Let V be the set of interference vertices of G .

There exists an edge $e = (v, w)$ between $v, w \in V$ if and only if the contexts enclosing v and w respectively, are not *control-dependent* && not *data-dependent*.

Edge Weight

To each edge $e = (v, w) \in E$ a weight $A(e)$ is associated:

$$A(e) = EF \times DW(e)$$

EF – execution frequency
 $DW(e)$ – distance weight =

$$\begin{cases} 1 & \text{if } v \text{ and } w \text{ are in the different statements} \\ 2 & \text{if } v \text{ and } w \text{ are in the same statement} \end{cases}$$

Previous Work

Powell, Lee, Newman, 1992 (assembly code; alternating fashion)
Saght, Chow, Lee, 1996 (post-pass approach; greedy algorithm)
Sudarsanam, Maly, 2000 (with register allocation; simulated annealing)
Leupers, Kotte, 2001 (before register allocation; ILP)
Zhuang, Fande, Greenland, 2002 (post-pass approach; maximum spanning tree)
Cho, Park, Whalley, 2002 (with register allocation; maximum spanning tree and graph coloring)
Zhuge, Xiao, Sha, 2002 (model based on DFC; scheduling with variable repartition)

Basic Features

- ▶ The partitioner is constructed as a separate optimization phase operating on the *high-level intermediate representation*.
- ▶ The *algorithm* of the partitioner is *global*, applied across the basic blocks to all functions of the program.
- ▶ *Global* variables and *static local* variables are handled.
- ▶ Array variables are treated as monolithic entities.

Partitioning Scheme

The partitioner is based on the concept of the *Interference Graph*.

- ▶ Constructing the interference graph.
- ▶ Partitioning the interference graph.
- ▶ Annotating variables with the partitioning information.

Partitioning the interference graph is solved as the combinatorial optimization problem Max-Cut :

$$\text{maximize } \sum_{e \in \text{Cut}(S,S)} A(e)$$

where $\text{Cut}(S,S)$ is the set of edges with one endpoint in S and the other endpoint in \bar{S} .

Max-Cut

Sahn and Gonzales, 1976 – the first feasible solution.
 Performance guarantee : $0.5 \times \text{optimal-value}$.

Goemans and Williamson, 1994 – the best algorithm.
 Performance guarantee : $0.878 \times \text{optimal-value}$.

Bur, Montero, and Zhang, 2001 – specialized version of the Goemans-Williamson technique.
 Performance guarantee : $0.878 \times \text{optimal-value}$.

Max-Cut (NP-Complete Problem)

Let $V = \{v_1, v_2, \dots, v_n\}$ be the set of vertices of G , and a_{ij} denotes the weight of the edge $(v_i, v_j) \in E$.
 When introducing cut vectors $x \in \{-1, 1\}^n$ with $x_i = 1$ for $v_i \in S$
 $x_i = -1$ for $v_i \in \bar{S}$
 then the algebraic formulation for Max-Cut :

$$\text{maximize } \frac{1}{2} \sum_{1 \leq i < j \leq n} a_{ij} (1 - x_i x_j)$$

subject to $x_i \in \{-1, 1\}, i = 1, \dots, n$

Max-Cut (Algebraic Formulation)

Let $\{v_1, \dots, v_k\} \subseteq V$ represent different accesses to the same variable. Then :

$\{v_1, \dots, v_k\}$ is merged to v
 $(v_i, w) \in E$ is redirected to $(v, w), i = 1, \dots, k$
 $(w, v_i) \in E$ is redirected to $(w, v), i = 1, \dots, k$

For $e = (v_i, w) :$
 $A(e)$ is defined to $\text{Max}(A(e_i)) \times k$
 $e_i = (v_i, w), i = 1, \dots, k$

Graph Merging

Partition the set of vertices of $G = (V, E)$ into two disjoint sets :

$$S \subseteq V \text{ and } \bar{S} = V - S$$

that the sum of the weights of all edges $(v, w) \in E$ connecting $v \in S$ and $w \in \bar{S}$, is maximal.

Graph Partitioning

Let (S, \bar{S}) be the result of the partitioning.
 Then each vertex :

$v \in S$ is assigned to memory bank X
 $v \in \bar{S}$ is assigned to memory bank Y

Memory Bank Assigning

The partitioning technique was evaluated on the simulator of the xDSPcore. The metrics which the performance is measured – the number of *total cycles*, and the number of *memory conflicts*. We did experiments with :

- ▶ DSPstone benchmark suite
- ▶ FFT Filters
- ▶ GSM Audio Codec

Experimental Results

Kernel	X-Allocating			Partitioning		
	Cycl.	X	Y	Cycl.	X	Y
dot-product	625	200	0	525	100	100
convolution	625	200	0	525	100	100
matrix.mult.1	5368	2100	0	5368	1000	1100
matrix.mult.2	5014	2010	0	4993	1100	910
mat1x3	85	24	0	76	9	15
lms	219	95	0	188	48	47
fir2dim	963	304	0	963	144	160
biquad.m_sects	71	38	0	66	21	17

DSPstone Kernels

Total Number of Cycles		Kernel	X-Allocating	Partitioning
Cycl.	X			
525	200	dot-product	525	100
525	200	convolution	525	100
5368	2100	matrix.mult.1	5368	1100
5014	2010	matrix.mult.2	4993	910
85	24	mat1x3	76	15
188	95	lms	188	47
963	304	fir2dim	963	160
66	38	biquad.m_sects	66	17

DSPstone Kernels

$$\begin{aligned} & \text{maximize } \frac{1}{2} \sum_{1 \leq i < j \leq n} a_{ij}(1 - y_{ij}) \\ & \text{subject to } y_i \geq 0, y_{ij} = 1, i = 1, \dots, n \\ & \uparrow \\ & \text{maximize } C \bullet X \\ & \text{subject to } \text{diag}(X) = e, X \succeq 0 \end{aligned}$$

Max-Cut (Semidefinite Program)

Solution (y_1, \dots, y_n) represents n points on the surface of the unit sphere corresponding to the set of vertices V .
 Construct the hyperplane: $H(r) = \{y \in \mathbb{R}^n : r^T y = 0\}$
 using the random vector r uniformly distributed on the unit sphere.
 Partition the vertex set :
 $S = \{v_i \in V : y_i^T r \geq 0\}$
 $\bar{S} = \{v_i \in V : y_i^T r < 0\}$

Max-Cut (Goemans' Randomized Solution)

- ▶ Exact method
- ▶ Alternate method
- ▶ Approximating iterative method (greedy)
- ▶ Semidefinite programming relaxation
- ▶ Semidefinite rank-2 relaxation

Implementation of Partitioning

Conclusion

- ▶ The algorithm :
- ▶ attempts to maximize the benefit of dual memory banks.
- ▶ is based on the partitioning the graph whose nodes represent variables and edges denote potential parallel accesses.
- ▶ operates on the *high-level intermediate representation*;
- ▶ its frame is *global*, not limited to basic blocks.
- ▶ finds a quite satisfying memory assignment – *improvement of memory cycles : 20% – 50%*
- ▶ *improvement of total cycles : 5% – 20%*

Future Plans

- ▶ To investigate the impact of the partitioning on the scheduling.
- ▶ To explore the memory partitioning for DSP architectures with interleaved memory banks.

of the Interference Graph: 13

X-Allocating		Alternate Allocating	
Cl.	X	Y	Confl.
941	91046	0	11053
Alternate Allocating		Alternate Allocating	
Cycl.	X	Y	Confl.
146322	48479	39920	8823 (80%)
Partitioning			
execution frequency (runtime)		without execution frequency	
977	44909	42518	8785 (79%)
152171	29957	61936	7427 (67%)
Total Number of Cycles			
Alternate Allocating	Alternate Allocating	Partitioning	
173394	155145 (89.5%)	154762 (89.2%)	

Size of the Interference Graph: 45

X-Allocating		Alternate Allocating	
Cycl.	X	Y	Confl.
1083592	150484	0	50693
Partitioning		Partitioning	
execution frequency (runtime)		without execution frequency	
1071731	81818	68642	290 (0.6%)
1073385	90273	60187	389 (0.8%)
Total Number of Cycles			
X-Allocating	Alternate Allocating	Partitioning	
1134285	1079968 (95.2%)	1072021 (94.5%)	

GSM Audio Codec

FFT Filter (256-point)

Executed Parallelism

Partitioning		X-Allocating	
total number of cycles	173394	154762 (89.2%)	87427
number of RW pairs	91046	22731 (52.0%)	14320 (32.8%)
number of parallel RW pairs	0	14320 (32.8%)	0
total number of cycles	1134285	1072021 (94.5%)	150460
total number of RW pairs	150484	50693 (67.1%)	58665 (78.0%)
number of parallel RW pairs	50693 (67.1%)	58375 (77.6%)	0