



H/S Codesign: Performance Analysis of Software

© Copyright
by
Margarida Jacome

The University of Texas at Austin



Margarida Jacome - UT Austin - 19



Performance Analysis

Princeton Sharad Malik, Wayne Wolf, Andrew Wolfe
(Cinderella project)

- 1 Introduction
- 1 Microarchitectural Modeling
 - ◆ Cache
 - ◆ CPU Pipelining
- 1 Experimental Results



Margarida Jacome - UT Austin - 19



Worst Case Analysis

- 1 **Worst case execution time (WCET) of a process running on a given processor**
 - ◆ **In real time systems: designer must verify if the WCET satisfies the timing deadlines**
 - » Many real-time OSs rely on this information for process scheduling
 - ◆ **In embedded system design: the WCET is often required to decide on Hw/Sw partitioning**



Worst Case Analysis

- 1 **Determining the actual WCET of a program**
 - ➡ ***Impractical***: requires simulating all possible combinations of *input data values* and *initial system states*
- 1 **Alternative: obtaining an estimate of the actual WCET by a static analysis of the program**
 - ◆ **Estimated WCET: must be tight and conservative (informative upper bound)**



■■■■ Problem Formulation



Determining the the estimated WCET of a given program on a given hardware system, assuming *uninterrupted execution*

- ◆ hardware system: includes the microprocessors and the memory systems
- 1 Problem has two components
 - ◆ Program path analysis
 - ◆ Microarchitectural modeling

■■■■ Estimating WCET

- 1 The problem of finding a program's estimated WCET is in general undecidable
- 1 Conditions for this problem to be decidable
 - ◆ absence of recursive function calls
 - ◆ bounded loops

Program Path Analysis

- 1 In order to tighten the estimated WCET, it is necessary to remove logically unfeasible program paths
 - ◆ path information provided by programmer
 - ◆ scope of path information provided by designer may have direct impact on the tightness of the estimated WCET
- 1 Statically feasible paths and programmer annotations
 - ◆ loop bounds;
 - ◆ interactions among different program statements (e.g., mutual exclusion of two statements)

Estimating WCET

- 1 Approach: the problem of solving the estimated WCET is converted into a set of Integer linear programming (ILP) problems
- 1 Discussion assumes (for now) a simple (pessimistic) microarchitectural model
 - ◆ execution time of an instruction is constant
 - » every instruction fetch results in a cache miss

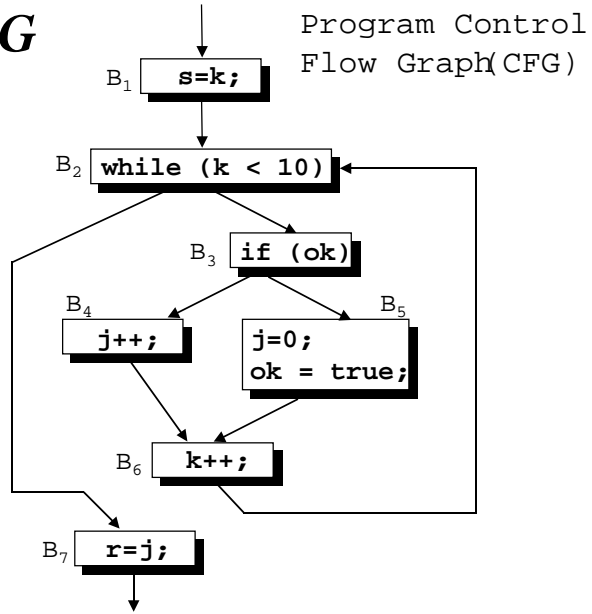
Program CFG

A basic block is a maximum contiguous sequence of instructions, with a single entry point at its start and a single exit point at its end.

```

/*k >= 0*/
s=k;
while (k < 10){
  if (ok)
    j++;
  else {
    j=0;
    ok = true;
  }
  k++;
}
r=j;

```



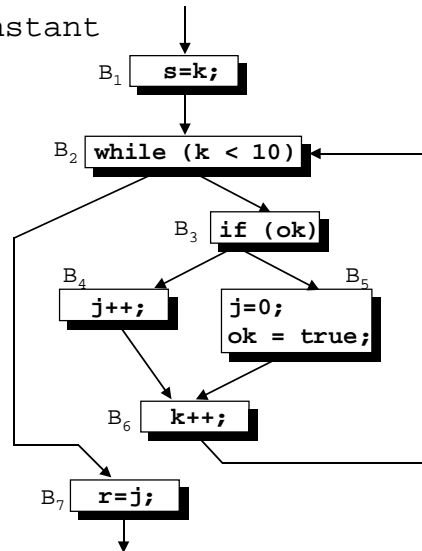
Estimating WCET: ILP Formulation

Assumption: each instruction takes constant time to execute

$$\text{Total execution time} = \sum_{i=1}^N c_i x_i$$

N = number of basic blocks

c_i = execution time of basic block i
 x_i = execution count of basic block i



Estimating WCET: ILP Formulation

Assumption each instruction takes constant time to execute

$$\text{Total execution time} = \sum_{i=1}^N c_i x_i$$

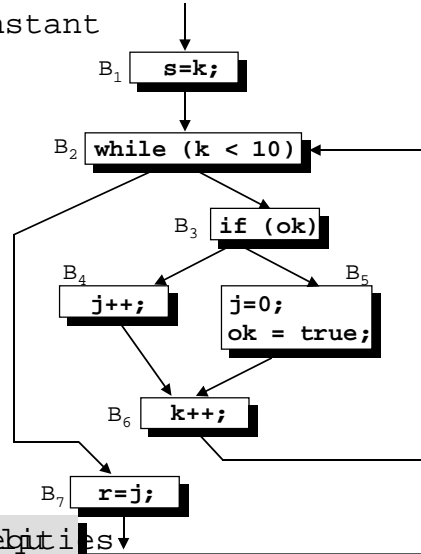
N = number of BBs

c_i = execution time of basic block i
 x_i = execution count of basic block i

Constrained by

- Program structure
- Values of program variables

constraints represented as linear equations



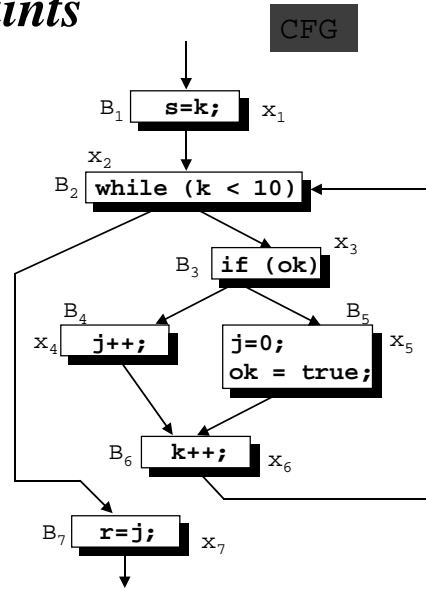
Execution Count of BBs

Two types of constraints are:

- 1 **Structural constraints**
 - ◆ derived from the program's control flow graph (CFG)
- 1 **Program functionality constraints**
 - ◆ provided by the user to specify loop bounds, and other path information

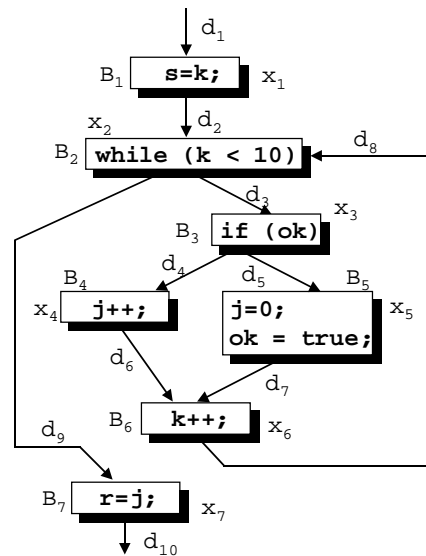
Structural Constraints

- Each node B_i defines an execution context



Structural Constraints

- Each node B_i defines an execution context
- Each edge defines a variable that counts the number of times the control flow passes through that edge

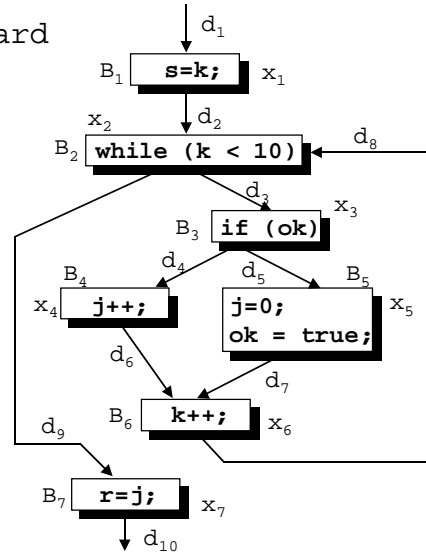


Structural Constraints

Analysis of CFG similar to standard network flow problem.

Execution count of node B_i is equal to:

- ⇒ Number of times that control enters the node (inflow)
- ⇒ Number of times that control exits the node (outflow)



Structural Constraints

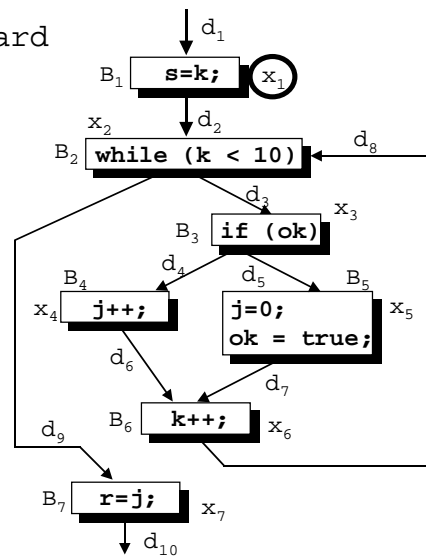
Analysis of CFG similar to standard network flow problem.

Execution count of node B_i is equal to:

- ⇒ Number of times that control enters the node (inflow)
- ⇒ Number of times that control exits the node (outflow)

$d_1 = 1$ - specifies that program is to be executed once

$$x_1 = d_1 = d_2$$

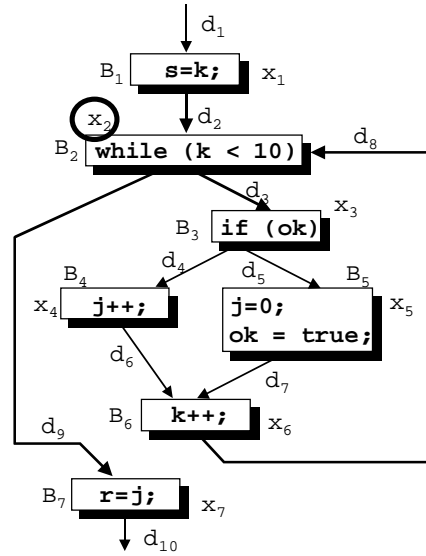


Structural Constraints

$$d_1 = 1$$

$$x_1 = d_1$$

$$x_2 = d_2 + d_8 = d_1 + d_8$$



Structural Constraints

$$d_1 = 1$$

$$x_1 = d_1$$

$$x_2 = d_2 + d_8 = d_1 + d_8$$

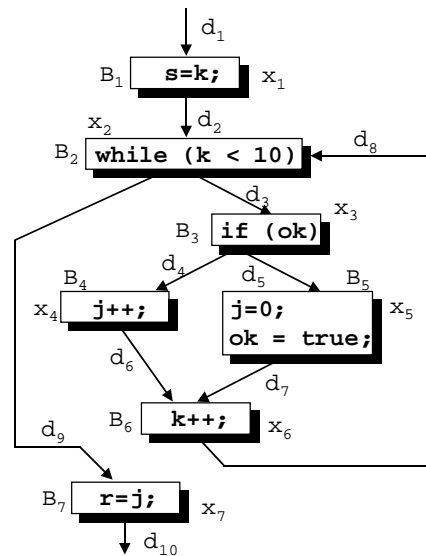
$$x_3 = d_3 + d_5 = d_1 + d_5$$

$$x_4 = d_4$$

$$x_5 = d_5$$

$$x_6 = d_6 + d_7 = d_4 + d_5$$

$$x_7 = d_9$$



Constraints Execution Count of BBs

- 1 Structural constraints
 - ◆ derived from the program's control flow graph (CFG)

- ➔ 1 Program functionality constraints
- ◆ provided by the user to specify loop bounds, and other path information

Program Functionality Constraints

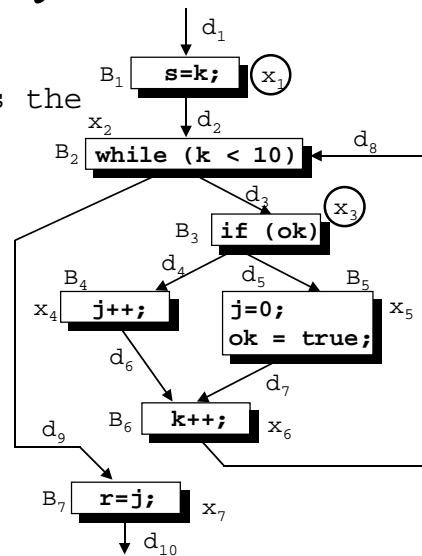
➔ Loop Bounds

Since k is positive before it enters the loop, the loop body will be executed between 1 and 10 times and the loop is exited.

$$0x_1 \leq x_3 \leq 10x_1$$

```

/*k >= 0*/
s=k;
while (k < 10){
  if (ok)
    j++;
  else {
    j=0;
    ok = true;
  }
  k++;
}
r=j;
    
```



Program Functionality Constraints

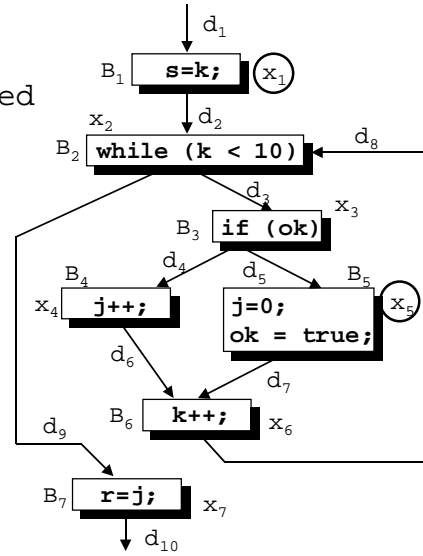
Other path information

Ex.1: The statement $s=k;$ can be executed at most once inside the loop

$$x_5 \leq 1 \quad x_3$$

```

/*k >= 0*/
s=k;
while (k < 10){
  if (ok)
    j++;
  else {
    j=0;
    ok = true;
  }
  k++;
}
r=j;
    
```



sti99 - 19

Program Functionality Constraints

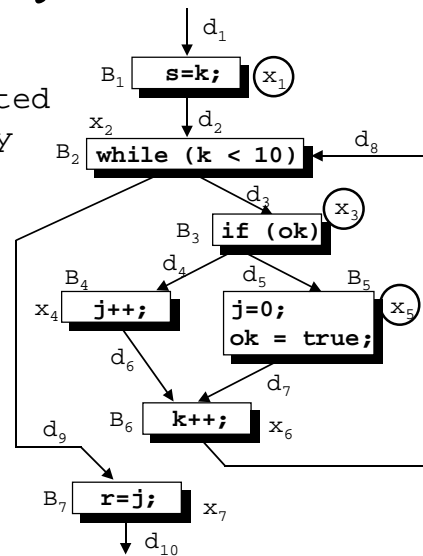
Other path information

Ex.2: If the else statement is executed the loop will be executed exactly 5 times

$$(x_5 = 0) \vee x_5 \geq 1 \wedge x_3 = 5 \wedge x_1$$

Non linear constraint

disjunction of linear constraint sets
(at least one constraint must be satisfied)



■■■■ Solving the Constraints

- 1 Because of the AND and OR operators, the program functionality constraints may, in general, be a disjunction of conjunctive constraint sets
- 1 To solve the estimated WCET, each set of the functionality constraints is combined with the set of structural constraints
 - ◆ the combined set is passed to the ILP solver with (1) to be maximized

$$\text{Total execution time} = \sum_i^N c_i x_i \quad (1)$$

→ ILP solver returns max value of expression and the BB counts

Margarida Jacome - UT Austin - 19

■■■■ Solving the Constraints

- 1 Procedure is repeated for every functionality constraint set
- 1 The maximum over all these running times is the estimated WCET

$$\text{Total execution time} = \sum_i^N c_i x_i$$

Margarida Jacome - UT Austin - 19

Time Complexity

- 1 Total time required to solve the estimated WCET depends on the number of functionality constraint sets and the time to solve each constraint set..
- 1 Number of functionality constraint sets doubles for each functionality constraint with a disjunction operator

Microarchitectural Modeling

Objective: model CPU pipeline and cache memory systems and determine execution times of BBs

Assumption: a cache miss stalls the CPU pipeline

Allows to divide execution of a BB into parts:

- Time for cache penalty (program cache)
- Memory-invariant real execution time in CPU pipeline

■■■■ Directed-Mapped Instruction Cache Analysis

- 1 Previous cost function (total execution time) is modified
- 1 Linear constraints representing the cache memory behavior (cache constraints) are added

■■■■ Modified Cost Function

- 1 Without cache, the factors that affects total execution time are the execution counts of BBs and their execution times
- 1 With cache, the execution time of an instruction may vary, depending on whether it results in a cache hit or a cache miss
 - ◆ for each instruction, subdivide original execution counts into counts of the *numbers of cache hits and cache misses*
 - ◆ consider the cache hit and cache miss execution time of each instruction

Modified Cost Function

1 Instructions inside a BB

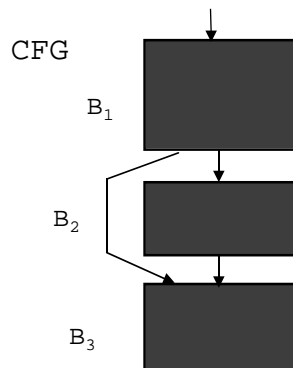
- ◆ have the same execution count
- ◆ may have different cache miss and cache hit counts
 - » they may map to different cache lines, each of which may have different cache activity



Different groups of instructions i
utilized for cache analysis is new way
of atomic structure

Line-Block (l-block)

- ### 1 Defines a contiguous sequence of instructions within the same BB that are mapped to the same line in the instruction cache
- ◆ the instructions within an *l-block* will always be executed the same number of times and have the same cache hit and cache miss counts

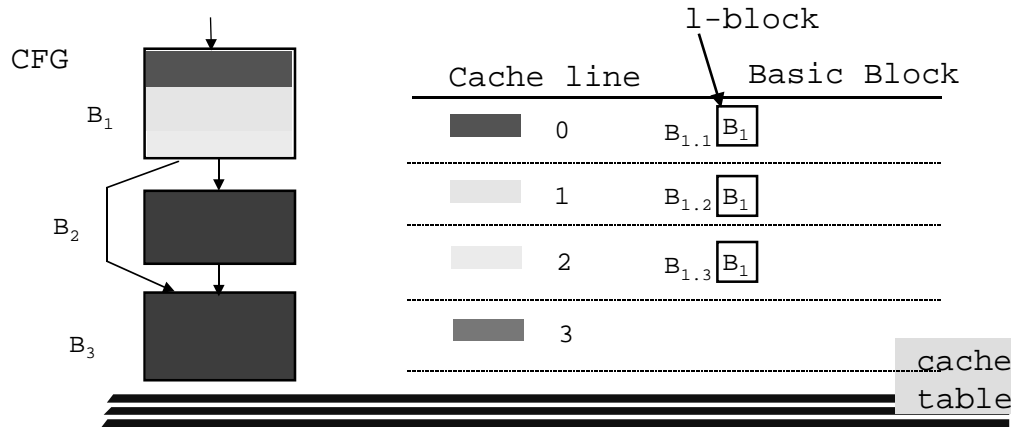


Example:

- ☛ CFG with 3 BBs
- ☛ Instruction cache with 4 lines

Identifying *l*-blocks

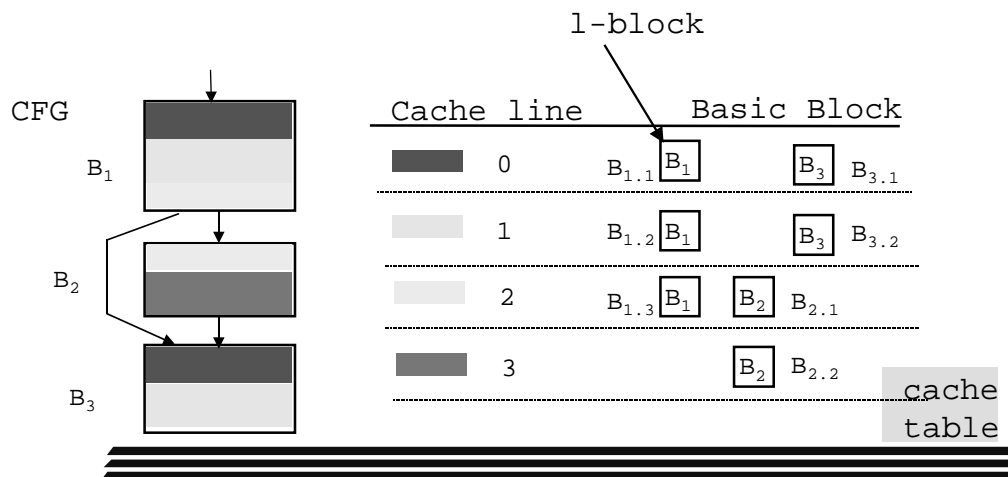
- For each basic block, find all the cache lines that instructions within it map to
 - ◆ Add an entry on these cache lines in the cache table



Margarida Jacome - UT Austin - 19

Line-Block (*l*-block)

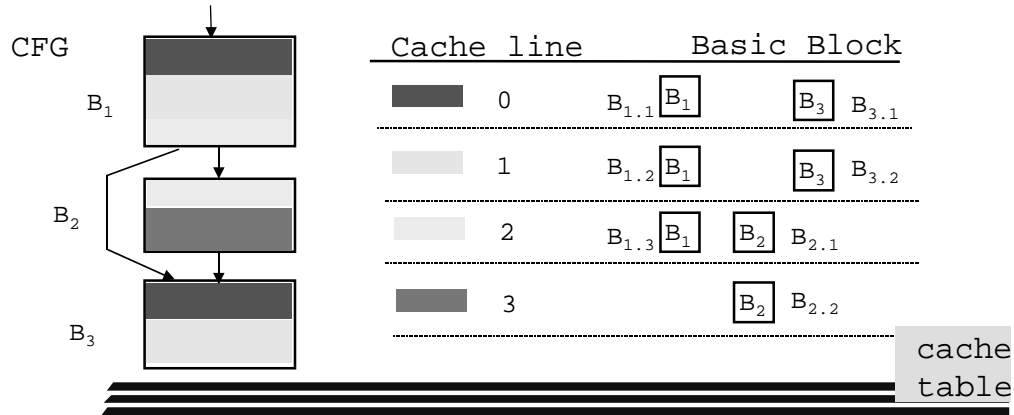
- Basic block B_1 is partitioned into three *l*-blocks: $B_{1.1}$, $B_{1.2}$ and $B_{1.3}$



Margarida Jacome - UT Austin - 19

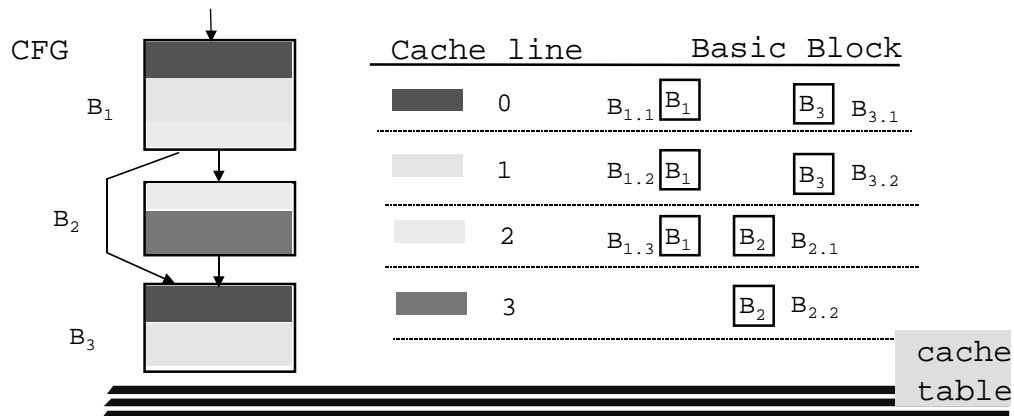
Line-Blocks (l-blocks)

- Two basic blocks are said to conflict with each other if the execution of one l-block will displace the other l-block in the instructions cache
 - for this to happen, they must map to the same cache line
 - Example: L-block $B_{1,1}$ conflicts with $B_{3,1}$



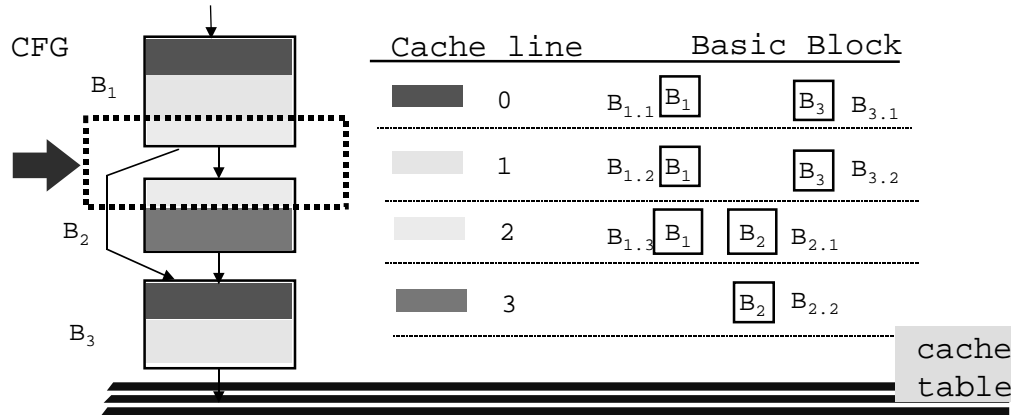
Line-Blocks (l-blocks)

- L-Block $B_{2,2}$ does not conflict with any other blocks (once it is in cache, it will never be displaced)



Line-Block (l-block)

- 1 L-blocks $B_{1,3}$ and $B_{2,1}$ occupy only a partial cache line
- ◆ A cache miss during the execution of either l-block will cause the system to load the instructions of both l-blocks into cache. Later executions will be cache hits.

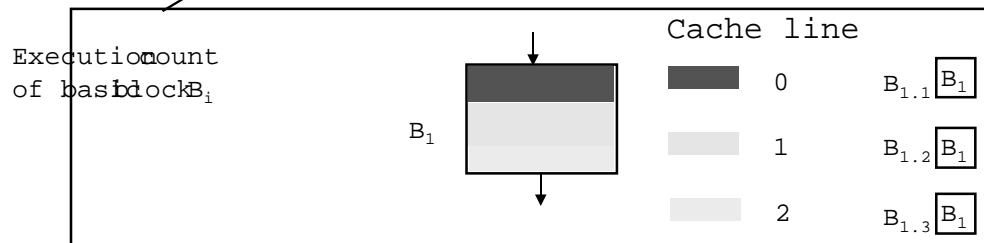


Margarida Jacome - UT Austin - 19

Execution Count of L-Blocks

- 1 Execution count of l-block $B_{i,j}$ is denoted as $x_{i,j}$
 - 1 Cache hit count is denoted $x_{i,j}^{hit}$
 - 1 Cache miss count is denoted $x_{i,j}^{miss}$
- (number of l-block of basic block B_i)

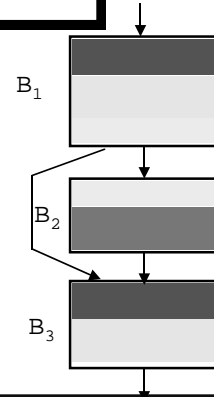
$$x_i = \sum_{j=1}^{n_i} x_{i,j} = \sum_{j=1}^{n_i} x_{i,j}^{hit} + \sum_{j=1}^{n_i} x_{i,j}^{miss}, \quad 1 \leq i \leq n$$



██████ Total Execution Time

$$\text{Total execution time} = \sum_i^N \sum_j^{n_i} (c_{i,j}^{\text{hit}} x_{i,j}^{\text{hit}} + c_{i,j}^{\text{miss}} x_{i,j}^{\text{miss}})$$

N = number of BBs
 n_i = number of blocks of BB B_i
 $x_{i,j}^{\text{hit}}$ = cache hit count of block $B_{i,j}$
 $x_{i,j}^{\text{miss}}$ = cache miss count of block $B_{i,j}$
 $c_{i,j}^{\text{hit}}$ = hit cost of block $B_{i,j}$
 $c_{i,j}^{\text{miss}}$ = miss cost of block $B_{i,j}$



(discussed later CPU pipeline)

██████ Constraints

$$\text{Total execution time} = \sum_i^N \sum_j^{n_i} (c_{i,j}^{\text{hit}} x_{i,j}^{\text{hit}} + c_{i,j}^{\text{miss}} x_{i,j}^{\text{miss}})$$

- 1 No changes to previous structural and program functionality constraints

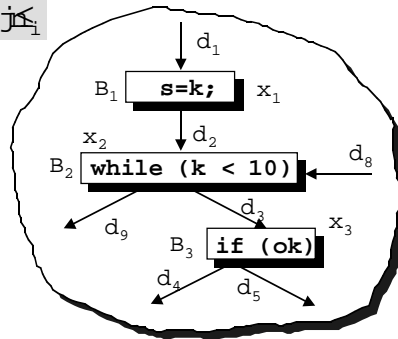
$$x_i = x_{i,j}^{\text{hit}} + x_{i,j}^{\text{miss}}, \quad 1 \leq j \leq n_i$$

Structural

$$\begin{aligned} d_1 &= 1 \\ x_1 &= d_1 \\ x_2 &= d_2 + d_8 \\ x_3 &= d_3 + d_9 \end{aligned}$$

Functional

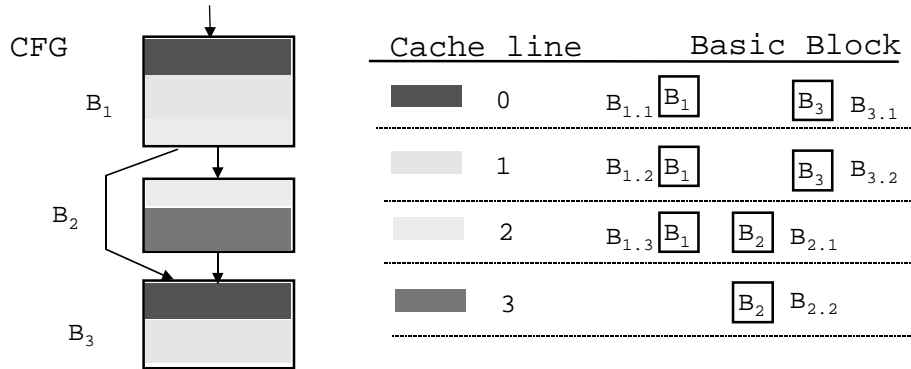
$$(x_5 = 0) \quad x_5 \geq 1 \quad \& \quad x_3 = 5 \quad x_1$$



Cache Constraints

$$\text{Total execution time} = \sum_i^N \sum_j^{n_i} (c_{i,j}^{\text{hit}} x_{i,j}^{\text{hit}} + c_{i,j}^{\text{miss}} x_{i,j}^{\text{miss}})$$

1 Additional "cache constraints" required...

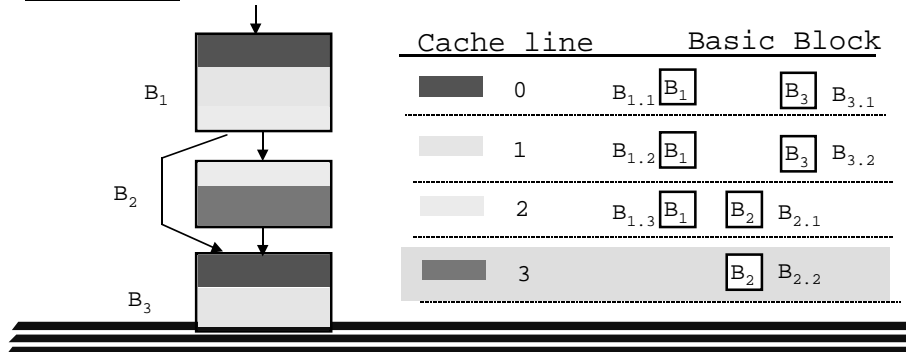


Cache Constraints

For each cache line l_i

- ➔ If there is only one block B_k that maps to it, the first execution of the block may cause a cache miss, and all subsequent executions will result in cache hits

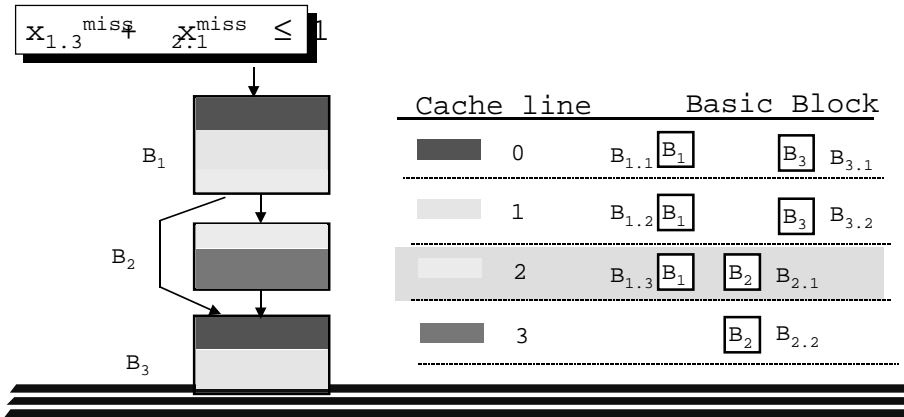
$$x_{k,l}^{\text{miss}} \leq 1 \quad \text{and} \quad x_i = x_{i,j}^{\text{hit}} + x_{i,j}^{\text{miss}}$$



Cache Constraints

For each cache line l_i

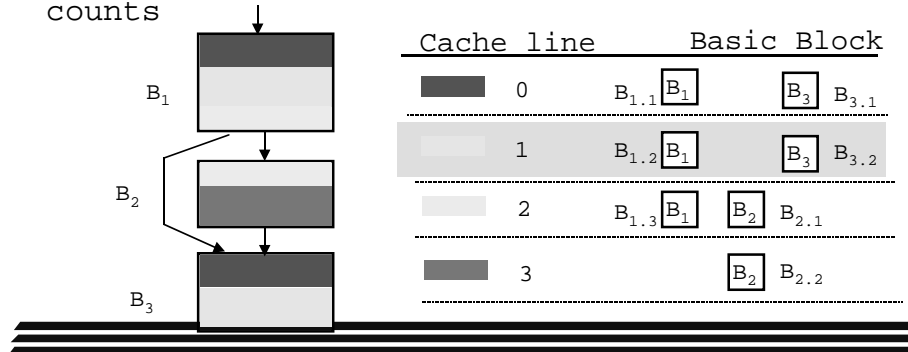
- ➔ If two or more non-conflicting blocks (B_1 and B_2) map to the execution of one of them will add the entire cache line count.



Cache Constraints

For each cache line l_i

- ➔ If two or more conflicting blocks (B_1 and B_2) map to it, the execution will be affected by the sequence in which these blocks are executed.
- ➔ The execution of other cache lines does not affect counts.



Cache Conflict Graph

1 A cache conflict graph is constructed for every cache line that contains two or more conflicting I-blocks

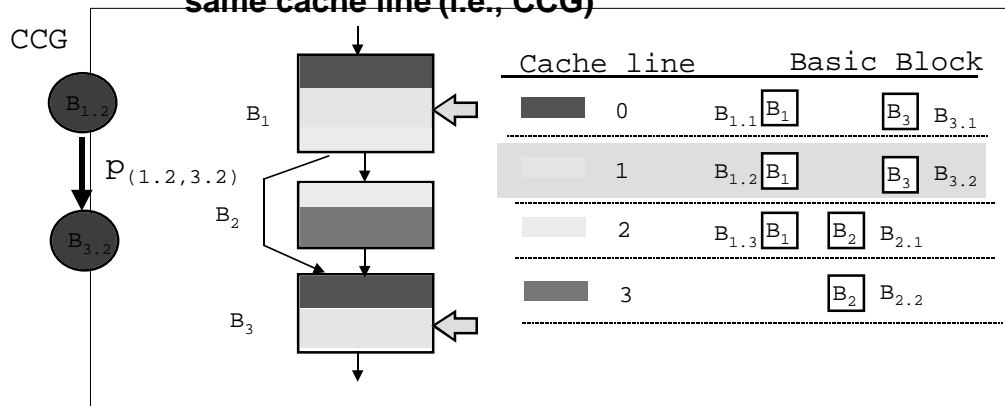
1 Nodes:

- ◆ a start node s and an end node e
 - » representing the start and the end of the program
- ◆ a node $B_{k,l}$ for every I-block that maps to that cache line

Cache Conflict Graph

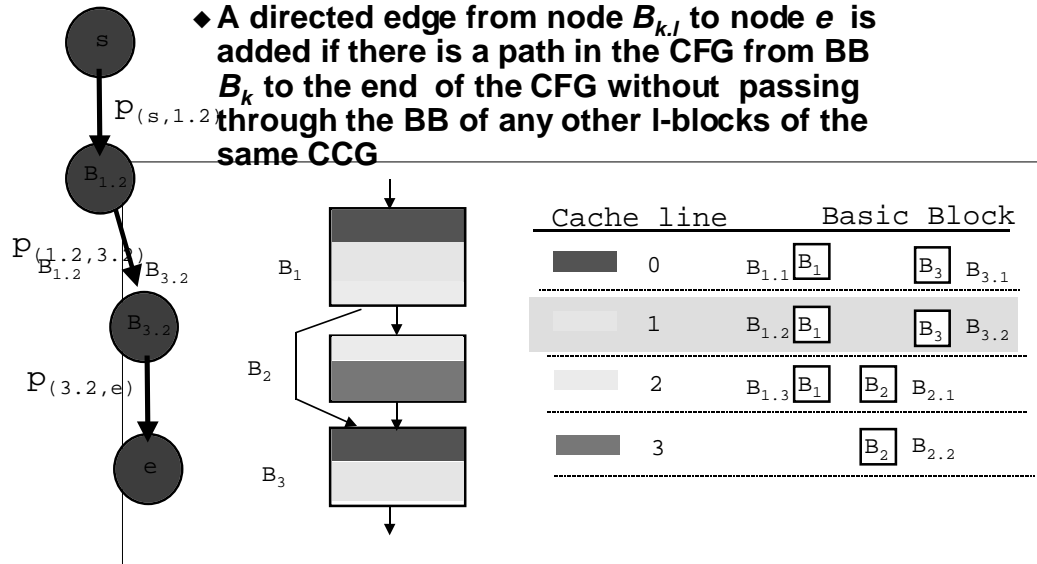
1 Edges

- ◆ A directed edge from node $B_{k,l}$ to node $B_{m,n}$ is added if there is a path in the CFG from BB B_k to BB B_m without passing through the BB of any other I-blocks of the same cache line (i.e., CCG)



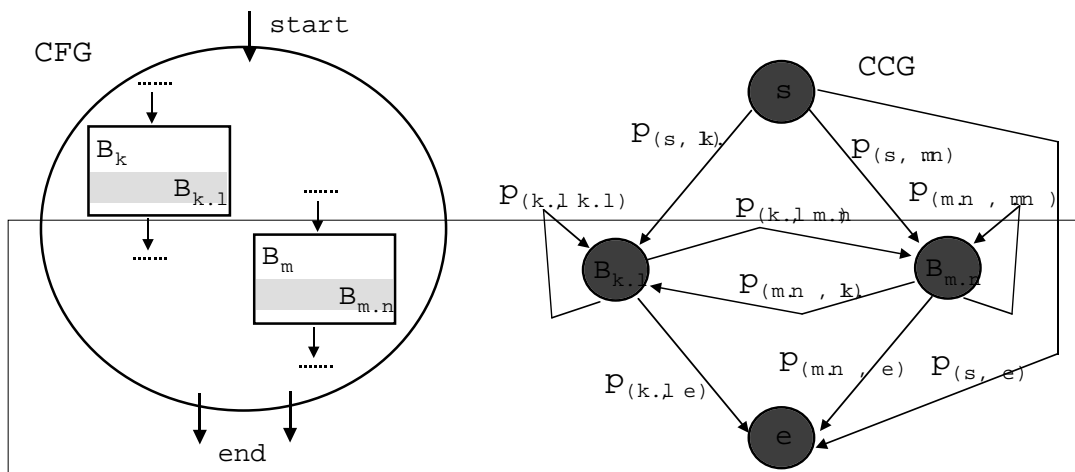
Cache Conflict Graph

1 Edges



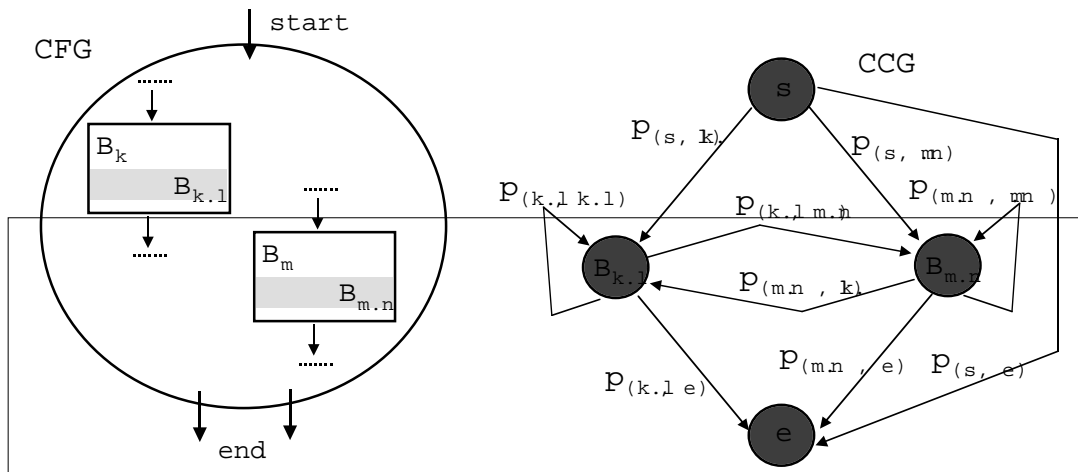
Cache Conflict Graph

➡ Possible CCG for two conflicting blocks



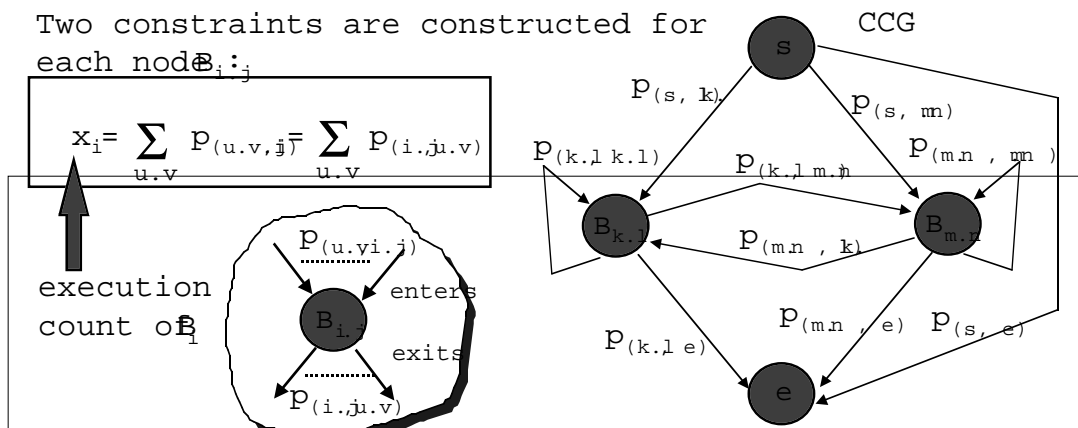
Cache Conflict Graph

- For each edge e , a variable p_e is assigned, representing the number of times that the control flow passes through that edge
- CCG has the structure of a network graph



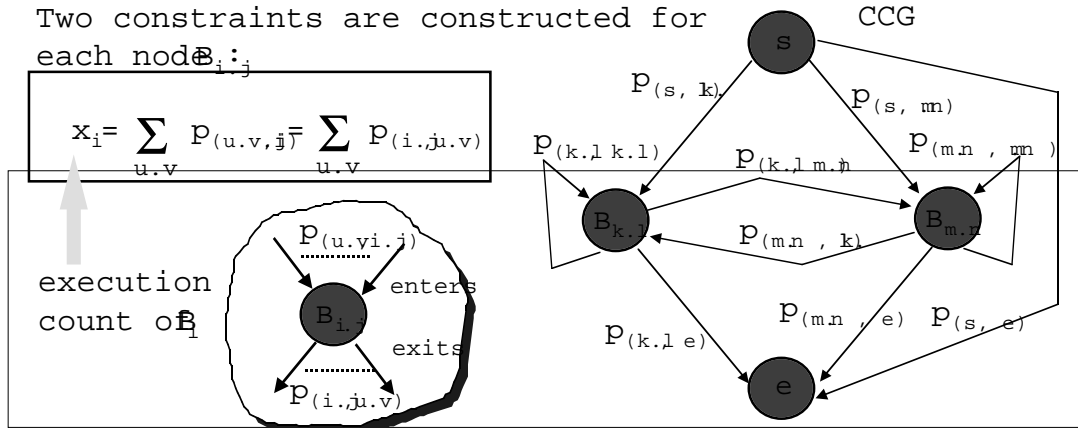
Cache Conflict Graph

- For each edge e , a variable p_e is assigned, representing the number of times that the control flow passes through that edge
- CCG has the structure of a network graph



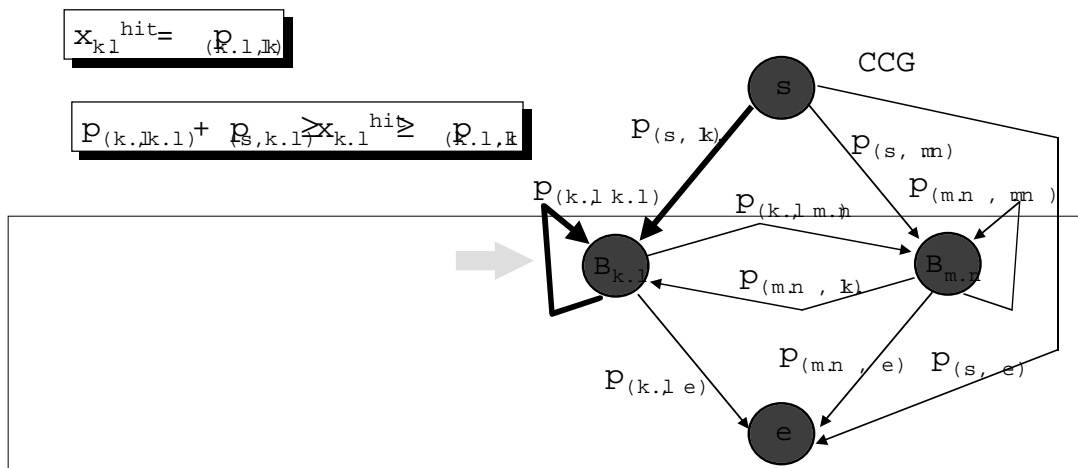
Cache Conflict Graph

- Constraints define the relationships between the program structural and program functionality constraints (via variables)



Cache Conflict Graph

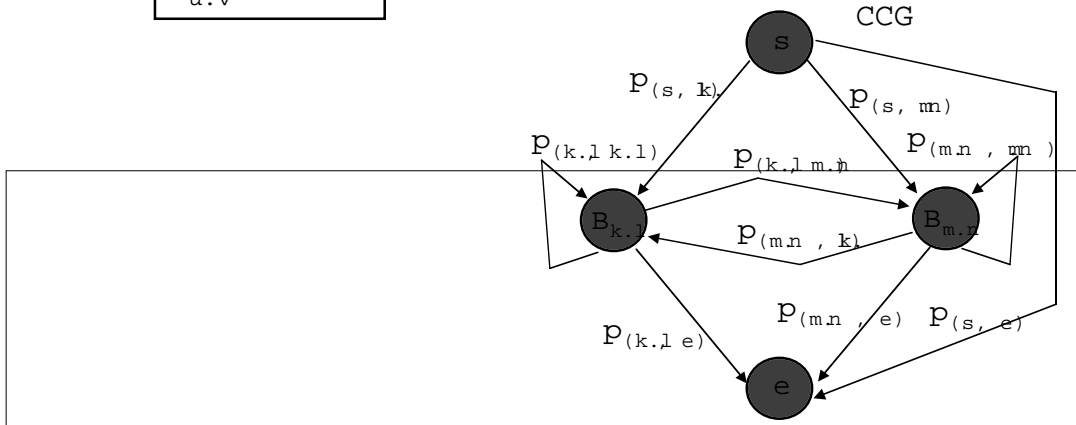
- The p-variable $p_{(k,l,k,l)}$ represents the number of times that the control flows into block $B_{k,l}$ after executing block $B_{k,l}$, without entering any other conflicting block in between



Cache Conflict Graph

Starting from indicating that the program is executed once

$$\sum_{u,v} p_{(s,u),v} = 1$$



Cache Constraints

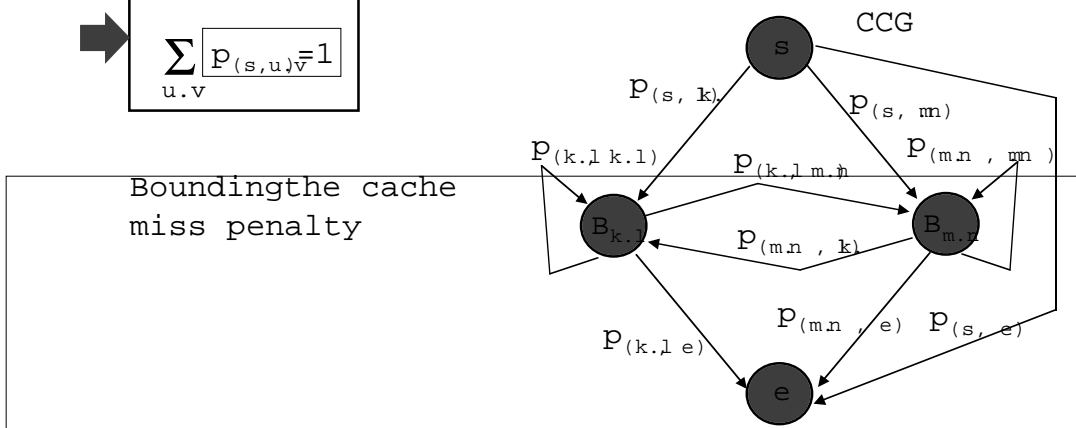
$$x_i = \sum_{u,v} p_{(u,v),j} \sum_{u,v} p_{(i,j),v}$$

$$x_{k,l}^{hit} = p_{(k,l),k}$$

$$p_{(k,l),k,l} + p_{(s,k,l)} x_{k,l}^{hit} \geq p_{(k,l),k}$$

$$\sum_{u,v} p_{(s,u),v} = 1$$

Bounding the cache miss penalty



Additional Constraints

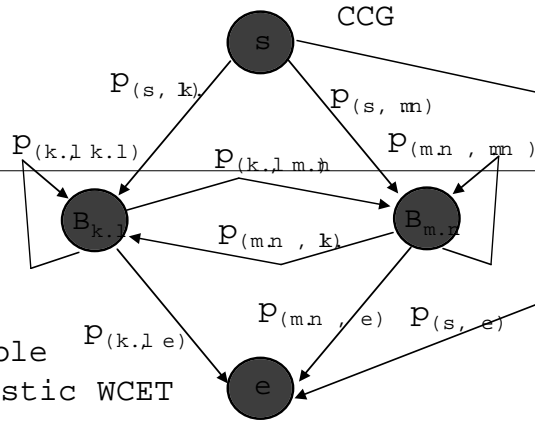
→ p-variables represent whether or there is a path between two conflicting blocks

The path represented by the p-variable may actually pass through a sequence of blocks...

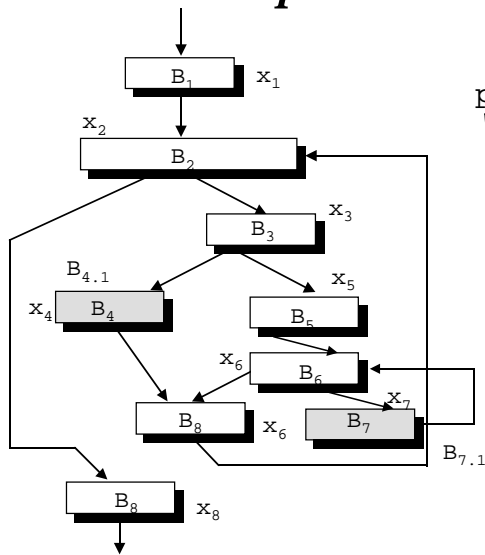
$$0 \leq p_{(m,n,k,l)} \leq \min(x_m, x_k)$$

Maximum value of the p-variable is bounded by the minimum of these BBs execution counts

Or ILP may return unfeasible 1-block count and pessimistic WCET



Example



$$x_4 = p(s,4.1) + p(4.1,4.1) + p(7.1,4.1) \\ = p(4.1,e) + p(4.1,1) + p(4.1,7)$$

$$x_7 = p(s,7.1) + p(7.1,7.1) + p(4.1,7.1) \\ = p(7.1,e) + p(7.1,1) + p(7.1,4)$$

Example

Cache Constraints

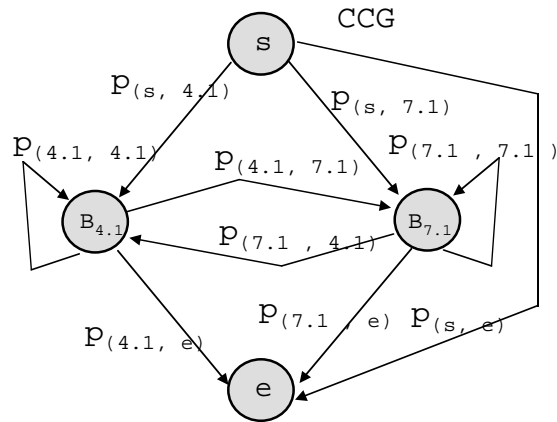
$$X_4 = p(s, 4, 1) + p(4, 1, 4) + p(7, 1, 4) \\ = p(4, 1, e) + p(4, 1, 1) + p(4, 1, 7)$$

$$X_7 = p(s, 7, 1) + p(7, 1, 7) + p(4, 1, 7) \\ = p(7, 1, e) + p(7, 1, 1) + p(7, 1, 4)$$

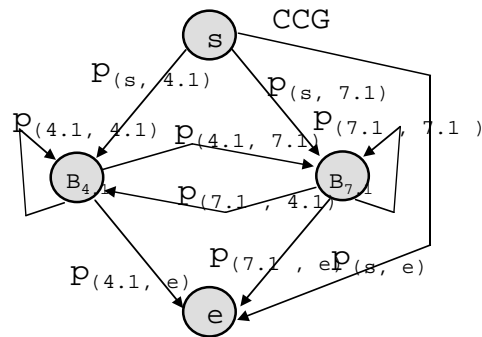
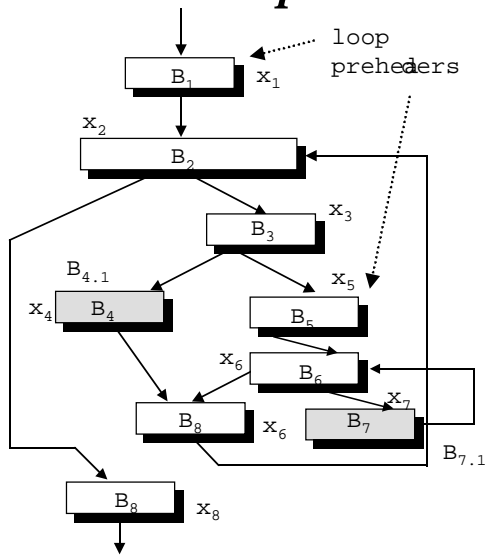
$$p(s, 4, 1) + p(s, 7, 1) + p(s, e) = 1$$

$$p(4, 1, 4) \leq x_4^{hit} \leq p(s, 4, 1) + p(4, 1, 4)$$

$$p(7, 1, 7) \leq x_7^{hit} \leq p(s, 7, 1) + p(7, 1, 7)$$



Example



Functional Constraints:

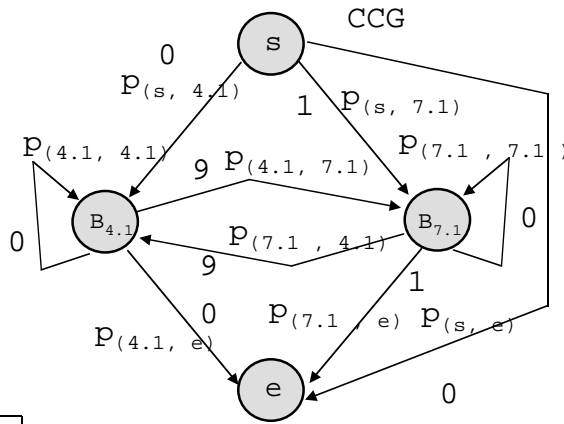
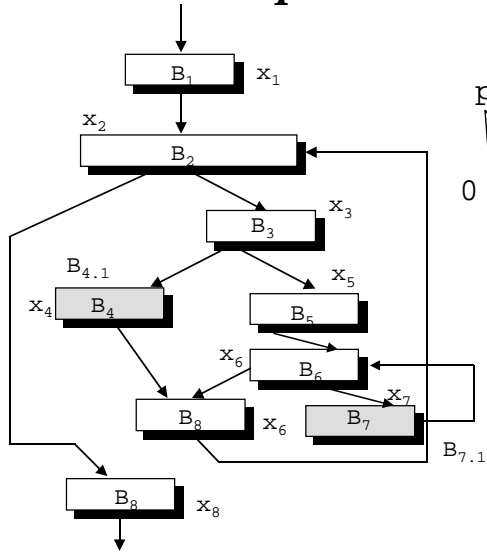
- Each loop will be executed 10 times each time they are entered
- Basic block B4 will be executed 9 times each time the outer loop is entered

$$x_3 = 10x_1 \quad x_7 = 10x_5$$

$$x_4 = 9x_1$$

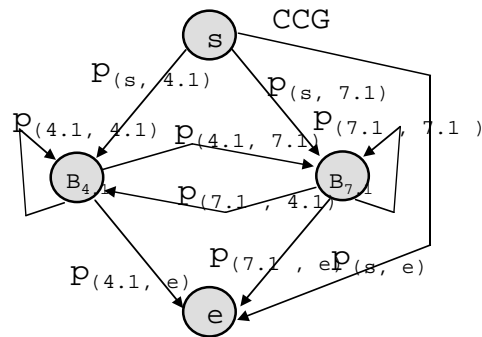
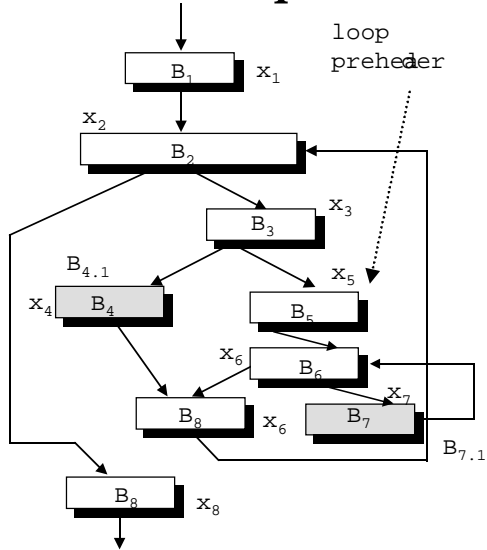


Example



WC Solution: maximum number of cache misses: impossible execution sequence

Example



Placing bounds on variables

- All paths enter the loop at the preheader
- Sum of the flows at nodes equal to exit count

$$p(7.1, 4.1) + p(4.1, 7.1) \leq 1$$

CPU Pipelining

- 1 **Assumption: time required to execute a sequence of instructions in the CPU pipeline is always a constant through the execution of the program**

$$\text{Total execution time} = \sum_i^N \sum_j^{n_i} (c_{i,j}^{\text{hit}} x_{i,j}^{\text{hit}} + c_{i,j}^{\text{miss}} x_{i,j}^{\text{miss}})$$

N = number of BBs
 n_i = number of blocks of BB B_i
 $x_{i,j}^{\text{hit}}$ = cache hit count of block $B_{i,j}$
 $x_{i,j}^{\text{miss}}$ = cache miss count of block $B_{i,j}$
 $c_{i,j}^{\text{hit}}$ = hit cost of block $B_{i,j}$
 $c_{i,j}^{\text{miss}}$ = miss cost of block $B_{i,j}$

CPU Pipelining

$c_{i,j}^{\text{hit}}$ = hit cost of block $B_{i,j}$
 $c_{i,j}^{\text{miss}}$ = miss cost of block $B_{i,j}$

- ✦ Found by adding up the effective execution times of the instructions in the block
- ✦ Some instructions (especially floating point instructions) are dependent
- ★ conservative approach: worst case effective time max (worst case execution time of floating point operation maybe 30% its average execution time).
- ✦ Additional time is added to the last block of each BB as to ensure that the buffered load/store instructions are completed when the control reaches the end of the BB

■■■■ CPU Pipelining

$$\begin{aligned} c_{ij}^{\text{hit}} &= \text{cost of } l\text{-block} \\ \rightarrow c_{ij}^{\text{miss}} &= \text{miss cost of } l\text{-B} \end{aligned}$$

- Equal to corresponding hit cost + time needed to load instructions of the block(s) in the cache memory

■■■■ Cinderella

- 1 Estimates the WCET of programs running on an Intel QT960 development board containing
 - ◆ 20 MHz Intel i960
 - » 32-bit RISC processor
 - » on-chip 512-byte direct-mapped instruction cache, organized as 32x16-byte lines
 - » 4 stage instruction pipeline
 - » floating point unit
 - ◆ 128K of main memory
- 1 Uses a public domain ILP solver

Experimental Results

- 1 “Actual WCET” was identified by authors, and the program execution time was measured for this worst case data set

Margarida Jacome - UT Austin - 19

Benchmarks

Function	Description	Lines	Bytes
check_data	Example from Pak's Thesis	17	88
pickrpt	Insertion sort	15	80
line	Line drawing routine Gupta's Thesis	143	1,556
circle	Circle drawing routine Gupta's Thesis	88	1,588
fft	Fast Fourier Transform	56	544
des	Data Encryption Standard	185	1,796
fullsearch	MPEG2 encoder frame search routine	204	1,436
whetstone	Whetstone benchmark	245	2,760
drystone	Drystone benchmark	480	1,360
matgen	Matrix routine in Linpac benchmark	50	248

Margarida Jacome - UT Austin - 19



Experimental Results

tighter bounds

Function	Measured WCET	Estimated WCET with cache analysis	Estimated WCET without cache analysis
check_data	4.41×10^2	4.91×10^2	11.9×10^2
picksrt	1.79×10^3	1.82×10^3	5.01×10^3
line	4.85×10^3	6.09×10^3	9.15×10^3
circle	1.45×10^4	1.53×10^4	1.59×10^4
fft	2.08×10^6	2.71×10^6	4.04×10^6
des	2.42×10^5	3.66×10^5	6.69×10^5
fullsearch	6.25×10^4	9.57×10^4	29.0×10^4
whetstone	6.83×10^6	10.2×10^6	14.9×10^6
drystone	5.52×10^5	7.53×10^5	13.3×10^5
matgen	9.28×10^3	10.9×10^3	17.2×10^3

clock cycle



Experimental Results

small integer programs

Function	Measured WCET	Estimated WCET with cache analysis	Estimated WCET without cache analysis
check_data	4.41×10^2	4.91×10^2	11.9×10^2
picksrt	1.79×10^3	1.82×10^3	5.01×10^3
line	4.85×10^3	6.09×10^3	9.15×10^3
circle	1.45×10^4	1.53×10^4	1.59×10^4
fft	2.08×10^6	2.71×10^6	4.04×10^6
des	2.42×10^5	3.66×10^5	6.69×10^5
fullsearch	6.25×10^4	9.57×10^4	29.0×10^4
whetstone	6.83×10^6	10.2×10^6	14.9×10^6
drystone	5.52×10^5	7.53×10^5	13.3×10^5
matgen	9.28×10^3	10.9×10^3	17.2×10^3