



# Efficient Microarchitecture Modeling and Path Analysis for Real-Time Software

Yau-Tsun Steven Li

Sharad Malik

Andrew Wolfe

# Introduction

- Paper examines the problem of determining the bound on the worst case execution time (WCET) of a given program on a given processor.
- Two important issues in solving this problem:
  - Program path analysis
  - Microarchitecture modelling
- Method which address both issues is proposed

# Two issues involved in solving this problem.

## Program path analysis.

- This determines what sequence of instructions will be executed in the worst case scenario.
- Infeasible program paths removed from the solution search space
- done by a data flow analysis of the program
- analysis should provide a mechanism for program path annotations.

## Microarchitecture modeling

- Models the hardware system and computes the WCET of a given sequence of instructions
- becoming difficult to model
- most modern processors have pipelined instruction execution units and cached memory systems.

# Proposed Solution

- Address both issues
- determine a tight bound on a program's worst case execution time.
- Explicit path enumeration not necessary to obtain tight estimated WCET
- Method determine worst case execution count of instruction and from these counts computes the estimated WCET
- includes a direct-mapped instruction cache analysis
- uses an integer linear programming formulation to solve the problem.
- allows the user to provide program path annotations so that a tighter bound may be obtained.

# Program path analysis problem handling

- Pessimistic approach:

Used simple microarchitecture model that assumes the execution time of an instruction to be a constant, i.e., every instruction fetch is assumed to result in a cache miss.

- method uses the counting approach to compute the estimated WCET.

- method converts the problem of solving the estimated WCET into a set of **integer linear programming (ILP) problems**

# ILP Formulation

- Assumption: Each instruction takes a constant time to execute
- Instructions within a basic block are always executed together, their execution counts are always the same.
  - let  $x_i$  be the execution count of a basic block  $B_i$ , and  $C_i$  be the execution time of the basic block,
  - given that there are  $N$  basic blocks in the program,

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

- possible values of  $x_i$  are constrained by the program structure and the possible values of the program variables.
- If these constraints represented as linear inequalities,  $\rightarrow$  the problem of finding the estimated WCET of a program is reduced to an integer linear programming (ILP) problem

# Linear Constraints

Divided into two parts:

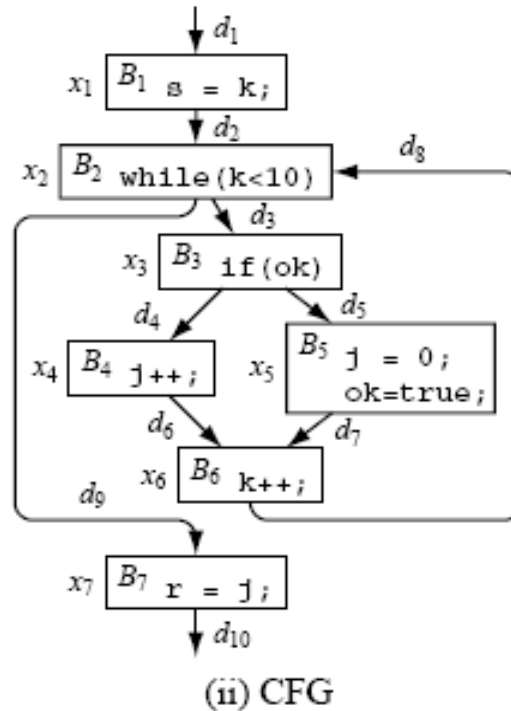
- **Program structural constraints,**
  - Derived automatically from the program's control flow graph (CFG)
- **program functionality constraints,**
  - provided by the user to specify loop bounds and other path information
  - or extracted from the program semantics.
- total time required to solve the estimated WCET depends on the number of functionality constraint sets and the time to solve each constraint set.
- the complexity of solving each ILP problem, an NP-hard problem.

# Example of Construction of these constraints

A conditional statement is nested inside a while loop

```

/* k>=0 */
S=k;
While (k < 10) {if (ok)
  J++;
  Else {
    J =0; ok=true;}
  K++;
} r=j;
    
```



*d<sub>i</sub>* = a count of the the number of times that the program control passes through that edge.

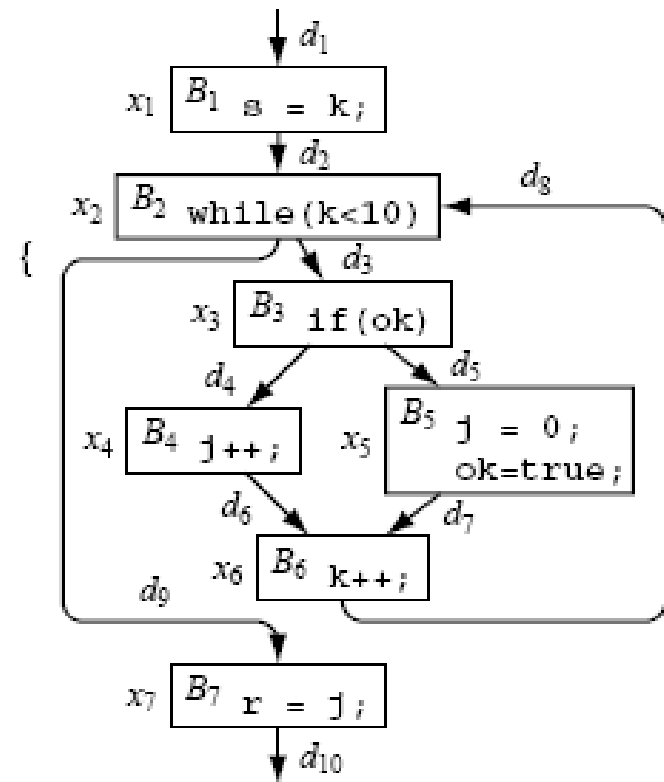
Each node in the CFG represents a basic block *B<sub>i</sub>*.  
 basic block execution count, *x<sub>i</sub>*, is associated with each node.



# Structural constraints

- Structural constraints can be derived from the CFG
- Fact: for each node  $B_i$ , its execution count is equal to the number of times that the control enters the node (inflow), and is also equal to the number of times that the control exits the node (outflow).
- structural constraints of this example
- Code fragment executed once, so  $d_1=1$

$$\begin{aligned}d_1 &= 1 \\x_1 &= d_1 = d_2 \\x_2 &= d_2 + d_8 = d_3 + d_9 \\x_3 &= d_3 = d_4 + d_5 \\x_4 &= d_4 = d_6 \\x_5 &= d_5 = d_7 \\x_6 &= d_6 + d_7 = d_8 \\x_7 &= d_9 = d_{10}.\end{aligned}$$



■ structural constraints do not provide any loop bound information

# Functional constraints

- Loop bound information can be provided by the user as a functionality constraint.

Example: since  $k$  is positive before it enters the loop, the loop body will be executed between 0 and 10 times each time the loop is entered


The constraints to specify this information are:  $0x_1 \leq x_3 \leq 10x_1,$

- The functionality constraints can also be used to specify other path information.

Example: the else statement ( $B5$ ) can be executed at most once inside the loop.

This information can be specified as:

$$x_5 \leq 1x_1.$$

- 
- To solve the estimated WCET, each set of the functionality constraint sets is combined (the conjunction taken) with the set of structural constraints.
  - The combined set is passed to the ILP solver with cost function to be maximized.

# Microarchitecture Modeling

- Previously, the modeling was simple because the execution time of an instruction was largely independent of others
- goal is to model the CPU pipeline and the cache memory systems and find out the execution times ( $C_i$ ) of the basic Blocks
- Method limited to model a direct-mapped instruction cache.
- can be extended to handle set associative instruction cache memory.

# Direct-mapped Instruction Cache Analysis

- To incorporate cache memory analysis in ILP model
  - need to modify the cost function
  - add a list of linear constraints, denoted as cache constraints, representing the cache memory behavior
- Modified Cost Function
  - With cache memory execution time of an instruction will be different depending on whether it results in a cache hit or cache miss.
  - need to subdivide the original instruction counts into counts of cache hits and misses.
- If cache hit and miss count and hit and miss execution time of instruction determined then tighter bound on execution time of program is established

# New type of atomic structure *line-block* (*l-block*) for analysis

- Adjacent instructions can be grouped together
- *l-block* is defined as a contiguous sequence of instructions within the same basic block that are mapped to the same line in the instruction cache.
- a basic block  $B_i$  is partitioned into  $n_i$  *l-blocks*. We denote these *l-blocks* as  $B_{i.1}, B_{i.2}, \dots, B_{i.n_i}$ .
- All instructions within an *l-block* will always have the same cache hit/miss counts, and the same total execution counts
- The cache hit and the cache miss counts of *l-block*  $B_{i.j}$  are denoted as  $x_{hit\ i.j}$  and  $x_{miss\ i.j}$

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

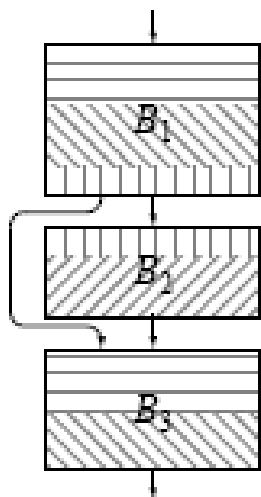
- the cache behavior can now be specified in terms of the new variables  $x_{hit\ i.j}$  and  $x_{miss\ i.j}$

New total execution time (Cost Function)

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

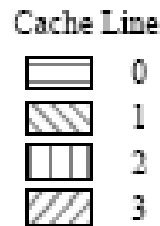
# Example showing how the I-blocks are constructed.

Each rectangle in the cache table represents a I-block.  
CFG with 3 basic blocks and instruction has 4 cache line



(i) CFG

B1.1  
B1.2  
B1.3  
B2.1  
B2.2  
B3.1  
B3.2



Cache line	Basic Block
0	$B_{1.1}$ $B_1$ $B_3$ $B_{3.1}$
1	$B_{1.2}$ $B_1$ $B_3$ $B_{3.2}$
2	$B_{1.3}$ $B_1$ $B_2$ $B_{2.1}$
3	$B_2$ $B_{2.2}$

(ii) Cache table

- any two I-blocks that map to the same cache line, they **conflict** with each other if the execution of one I-block will displace the cache content of the other.
- Otherwise, they are called **non-conflicting** I-blocks e.g. B1.3 and B2.1

# Cache Constraints

- used to constrain the hit/miss counts of the I-blocks.
- simple case : For each line only one I-block  $B_{k.l}$  mapping.  
First execution of this I-block may cause a cache miss and all subsequent executions will result in cache hits.

$$x_{k,l}^{miss} \leq 1.$$

- case : Two or more **non-conflicting** I-blocks map to the same cache line(  
 $B_{1.3}$  and  $B_{2.1}$

The execution of any of them will load all the I-blocks into the cache line. sum of their cache miss counts is at most one.

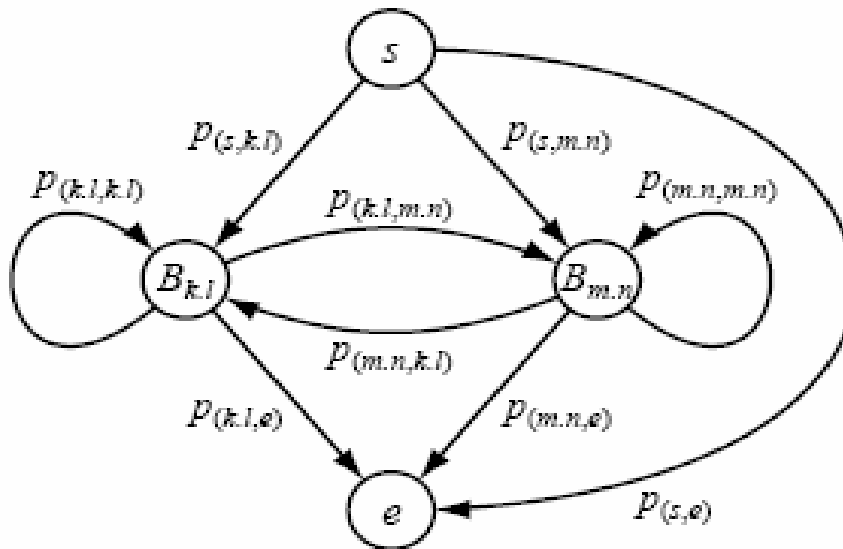
$$x_{1.3}^{miss} + x_{2.1}^{miss} \leq 1.$$

- Case: a cache line contains two or more **conflicting** I-blocks, the hit/miss counts of all the I-blocks mapped to this line will be affected by the sequence in which these I-blocks are executed.



## Cache Conflict Graph (Network flow graph)

- cache conflict graph (CCG) is constructed for every cache line containing two or more conflicting I-blocks. Example :Cache line contains 2 conflicting graph
- start node 's', an end node 'e', and a node  $B_{k,l}$  for every I-block  $B_{k,l}$  mapped to the same cache line.
- if there exists a path in the CFG from basic block  $B_k$  to basic block  $B_m$  without passing through the basic blocks of any other I-blocks of the same cache line



Program begins at S node. i) After executing other L-block from other cache line eventually reaches to one of conflicting graph

ii) After executing  $B_{k,l}$  may pass other I-block and reaches to  $B_{m,n}$  or directly passes to  $B_{m,n}$

$p(i, j, u, v)$  to count the number of times that the control passes through that edge

## continued

- At each node  $Bi. j$ , the sum of control flow going into the node must be equal to the sum of control flow leaving the node, and it must also be equal to the execution count of  $l$ -block  $Bi. j$ .
- two constraints are constructed at each node  $Bi. j$ :

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

This set of constraints is linked to structural and functionality constraints via the  $x$ -variables.

- Program executed once at start node

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

variable  $p(i. j, i. j)$  represents the number of times that the control flows into  $l$ -block  $Bi. j$  after executing  $l$ -block  $Bi. j$

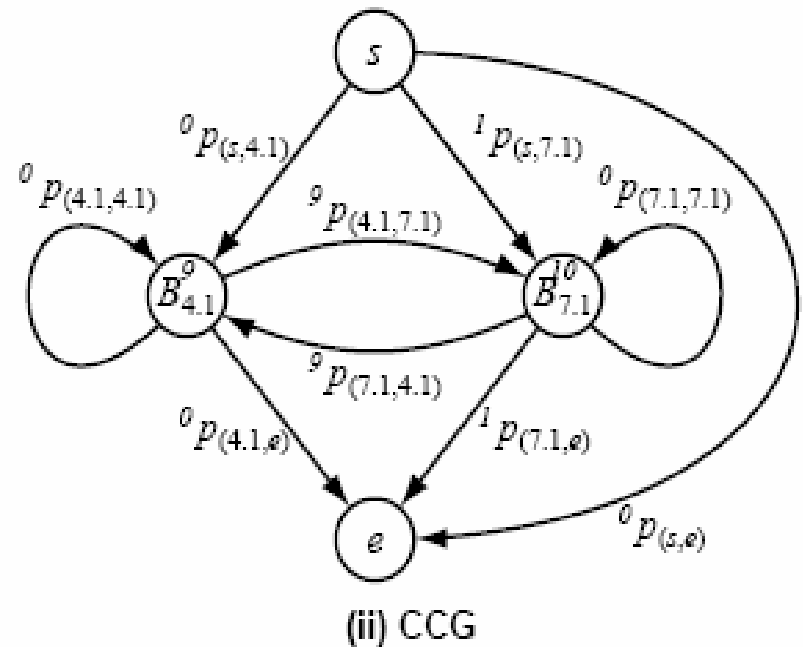
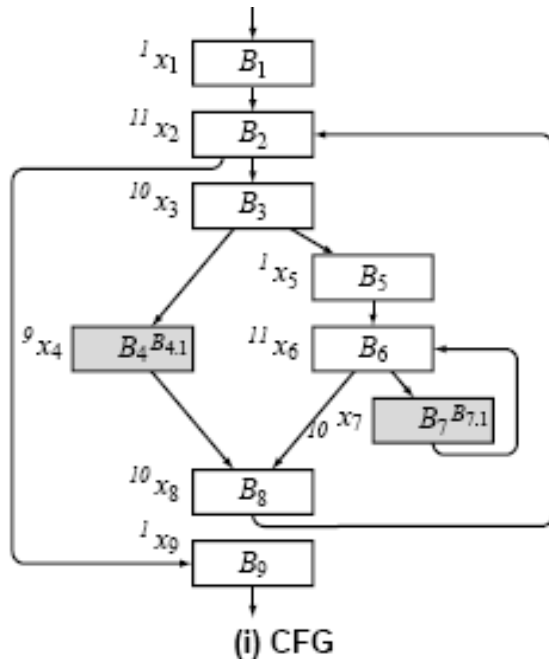
- If both edges (Bi.j, e) and (s, Bi.j) exists then the program variable p (s,i. j) *may* also be counted as a cache hit.

$$P(i,j,i,j) \leq x_{i,j}^{hit} \leq P(s,i,j) + P(i,j,i,j).$$

if any of edges (s,Bi. j) and (Bi. j ,e) does not exist, then  $x_{i,j}^{hit} = P(i,j,i,j).$

# Bounds on p-variables

- some path sequencing information can be expressed in terms of p-variables as extra functionality constraints
- Without the correct bounds, the solver may return an infeasible l-block count and an overly pessimistic estimated WCET.



example showing two conflicting l-blocks ( $B_{4.1}$  and  $B_{7.1}$ ) from two different loops.

The italicized numbers shown on the left of the variables are the pessimistic worst case solution returned from ILP solver.

For any variable  $p(i, j, u, v)$ , its bounds are:  $0 \leq p(i, j, u, v) \leq \min(x_i, x_u)$ .

A loop preheader is the basic block just before entering the loop. For instance, in the example shown in Fig. 4, basic block B1 is the loop preheader of the outer loop and basic block B5 is the loop preheader of the inner loop.

a constraint at loop preheader *B5* is needed

$$p(5, 7.1) + p(4.1, 7.1) \leq x_5.$$

# Interprocedural call

- function may be called many times from different locations of the program.
- Every function call is treated as if it is inlined.
- a function call is represented by an f - edge pointing to an instance of the callee function's CFG.
- edge has a variable 'fk ' which represents the number of times that the particular instance of the callee function is called

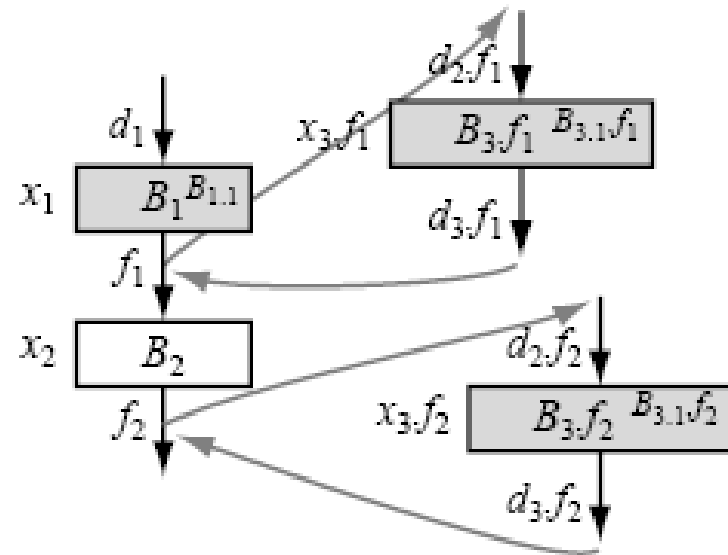
# Function inc is called twice in the main function

```

void main()
{
  B1  inc(&i);
  B2  inc(&j);
}

void inc(int *pi)
{
  B3  *pi++;
}

```

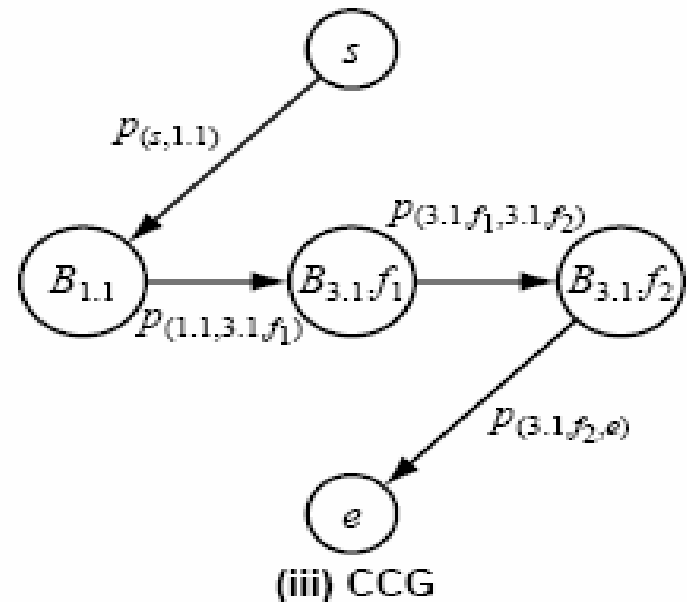
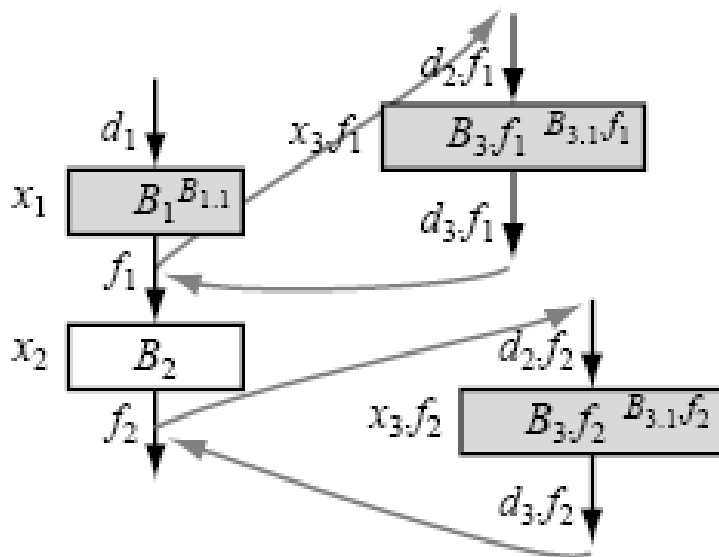


last equation above links the total execution counts of basic block  $B_3$  with its counts from two instances of the function

$$\begin{aligned}
 d_1 &= 1 \\
 x_1 &= d_1 = f_1 \\
 x_2 &= f_1 = f_2 \\
 d_2 \cdot f_1 &= f_1 \\
 x_3 \cdot f_1 &= d_2 \cdot f_1 = d_3 \cdot f_1 \\
 d_2 \cdot f_2 &= f_2 \\
 x_3 \cdot f_2 &= d_2 \cdot f_2 = d_3 \cdot f_2 \\
 x_3 &= x_3 \cdot f_1 + x_3 \cdot f_2
 \end{aligned}$$

# CCG

- CCG is constructed as before by treating each instance of I-block  $B_{i,j} \cdot f_k$  as different from other instances of the same I-block.
- In the example, if I-block  $B_{1,1}$  conflicts with I-block  $B_{3,1}$ , then since I-block  $B_{3,1}$  has two instances ( $B_{3,1} \cdot f_1$  and  $B_{3,1} \cdot f_2$ ), there will be 5 nodes in the CCG



Inc



# Cache constraints

- cache constraints and the bounds on  $p$  variables are constructed as before,
- the hit constraints are modified slightly. Edge going from  $B_{i,j} \cdot f_k$  to  $B_{i,j} \cdot f_1$  counted as cache hit of block  $B_{i,j}$
- The complete cache constraints derived from the example's CCG are

$$x_1 = x_{1.1}^{hit} + x_{1.1}^{miss}$$

$$x_2 = x_{2.1}^{hit} + x_{2.1}^{miss}$$

$$x_3 = x_{3.1}^{hit} + x_{3.1}^{miss}$$

$$x_{2.1}^{miss} \leq 1$$

$$x_1 = P(s, 1.1) = P(1.1, 3.1, f_1)$$

$$x_3 \cdot f_1 = P(1.1, 3.1, f_1) = P(3.1, f_1, 3.1, f_2)$$

$$x_3 \cdot f_2 = P(3.1, f_1, 3.1, f_2) = P(3.1, f_2, e)$$

$$P(s, 1.1) = 1$$

$$x_{1.1}^{hit} = 0$$

$$x_{3.1}^{hit} = P(3.1, f_1, 3.1, f_2)$$

# CPU Pipeline

- The CPU pipeline is considered to be relatively easy to model because it is only effected by adjacent instructions.

- As  $C_{i,j}^{hit}$  and  $C_{i,j}^{miss}$  must be constants

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

- hit cost  $C_{i,j}^{hit}$  of a l-block  $B_{i,j}$  is determined by adding up the effective execution times of the instructions in the l-block
- the effective execution times of some instructions, especially the the floating point instructions, are data dependent, a conservative approach is taken by assuming the worst case effective execution time
- Additional time is also added to the last l-block of each basic block so as to ensure that all the buffered load/store instructions are completed when the control reaches the end of the basic block.
- miss cost  $C_{i,j}^{miss}$  of the l-block is equal to the time needed to load the instructions of the l-block into the cache memory and to execute them in the CPU.

# Implementation

- cache analysis method has been implemented in a tool called cinderella4, which estimates the WCET of programs running
- The tool reads the subject program's executable code and constructs the CFGs and the CCG
- outputs the annotation files in which the 'x 'and 'f ' are labeled along with the program's source code
- user is then asked to provide loop bounds
- estimated WCET can thus be computed
- user can provide additional path information, if available, to tighten this bound.

# Experiment

Table 2: Estimated WCETs of Benchmark programs. All values are in units of clock cycles.

Function	Measured WCET	Estimated WCET with cache analysis	Estimated WCET w/o cache analysis
check_data	$4.30 \times 10^2$	$4.91 \times 10^2$	$11.9 \times 10^2$
piksort	$1.71 \times 10^3$	$1.74 \times 10^3$	$5.86 \times 10^3$
sort	$9.99 \times 10^6$	$27.8 \times 10^6$	$50.2 \times 10^6$
matent	$2.20 \times 10^6$	$5.46 \times 10^6$	$8.17 \times 10^6$
matent2	$1.86 \times 10^6$	$2.11 \times 10^6$	$4.46 \times 10^6$
stats	$1.16 \times 10^6$	$2.21 \times 10^6$	$2.95 \times 10^6$
fft	$2.20 \times 10^6$	$2.63 \times 10^6$	$3.97 \times 10^6$
jpeg_fdct_islow	$9.05 \times 10^3$	$9.11 \times 10^3$	$16.7 \times 10^3$
line	$4.84 \times 10^3$	$6.09 \times 10^3$	$9.15 \times 10^3$
circle	$1.45 \times 10^4$	$1.54 \times 10^4$	$1.59 \times 10^4$
des	$2.44 \times 10^5$	$3.70 \times 10^5$	$6.72 \times 10^5$
whetstone	$6.94 \times 10^6$	$10.5 \times 10^6$	$14.9 \times 10^6$
dhry	$5.76 \times 10^5$	$7.57 \times 10^5$	$13.3 \times 10^5$

# Conclusion

- tight bound on a program's WCET is estimated.
- small amount of pessimism due to
  - (i) insufficient path information from the user  
so that some infeasible program paths are considered, => can be reduced by providing more path information
  - (i) inaccuracy in microarchitecture modeling  
affects the accuracy of the values of  $C_{hi}.j$  and  $C_{mi}.j$  => reduced by a more sophisticated hardware model