# WCET'2002

# 2<sup>nd</sup> International Workshop on Worst-Case Execution Time Analysis (Satellite Event to ECRTS'02)

Technical University of Vienna, Austria June 18, 2002



# Message from the Workshop Chair

Welcome to the 2<sup>nd</sup> international Workshop on Worst-Case Execution Time (WCET) Analysis, a satellite event of the Euromicro Conference on Real-Time Systems. The workshop was held in Vienna, Austria, on the 18th of June 2002. This is the second event in the series after the successful first meeting held in Delft.

The aim of the workshop is to provide a forum for discussing current trends and issues related to the timing analysis of Real-Time Systems with special emphasis on bridging the gap between industry and academia. The meeting encourages debate and interaction between participants through short presentations followed by active discussion. The program of the workshop presents contributions on the following areas of timing analysis:

- Within the context of high-level analysis techniques contributions address path analysis techniques and issues related to object oriented programming models.
- On low-level analysis techniques the focus is on modelling timing behaviour of processor features such as cache effects, branch prediction and speculative execution.
- The industrial view presents timing requirements in the aerospace industry and current models of analysis and current tool support.

I would like to express my congratulations to all participants, authors, reviewers, and the organisation of the workshop (special thanks to Antoine Colin and Stefan Petters) that have made this event a successful one.

Dr. Guillem Bernat. University of York. England, UK

### Table of contents

Session I: High level analysis

- A Prototype Tool for Flow Analysis of C Programs. Jan Gustafsson, Björn Lisper, Nerina Bernmudo, Christer Sandberg and Linus Sjöberg. Mälardalen University, Västerås, Sweden.
- A novel Gain Time Reclaiming Framework Integrating WCET Analysis for Object-Oriented Real-Time Systems. Erik Yu-Shing Hu, Andy Wellings and Guillem Bernat. University of York, United Kingdom.
- A Unified Flow Information Language for WCET Analysis, Andreas Ermedahl, Jakob Engblom, Friedhelm Stappert.

Session II: Tools

- WCET Estimation from Object Code Implemented in the PERF Environment. Douglas Renaux, João Goés and Robson Linhares. Laboratory of Embedded Systems Innovation and Technology, Brasil.
- Status of the BOUND-T WCET Tool. Niklas Holsti and Sami Saarinen. Space Systems Finland Ltd., Espoo, Finland

Session III: Industrial views

- You Can't Control what you Can't Measure, or Why it's Close to Impossible to Guarantee Real-Time Software Performance on a CPU with On-Chip Cache. Nat Hillary and Ken Madsen. Applied Microsystems Corp./Wind River Sysrems Inc.
- Worst Case Execution Time Prediction.
   Marc Langenbach, Christian Ferdinand and Reinhard Wilhelm
- The European Space Agency's Involvement and interest in WCET and Scheduling Analysis.

Morter Rytter Nielsen, Eric Conquet and Jean-Loup Terraillon. ESA, Nordwjik, Netherlands.

Session IV: Low Level Analysis

- Cache Modelling vs Static Cache Locking for Schedulability Analysis in Multitasking Real-Time Systems. Isabelle Puaut. IRISA, Rennes, France.
- A Framework to Model Branch Prediction for WCET Analysis. Tulika Mitra and Abhik Roychoudhury. National University of Singapore, Singapore.
- *Difficulties in computing the WCET for Processors with Speculative Execution.* Christine Rochange and Pascal Sainrat. Institut de Recherche en Informatique de Toulouse, France.

Session V: Issues in WCET Analysis

- Why You Can't Analyze RTOSs without Considering Applications and Vice Versa. Jörn Schneider. Saarland University, Saarbrücken, Germany.
- *How Much Worst Case is Needed in WCET Estimation?* Stefan Petters. University of York, United Kingdom.
- Is WCET Analysis a Non-Problem? Towards New Software and Hardware Architectures. Peter Puschner. University of Vienna, Austria.

# Session I: High level analysis

### Isabelle Puaut

# Presentations

The "high level analysis" session was made of three presentations. The first one introduced ongoing work on the automatic determination of flow facts (loop bounds, infeasible paths) of C programs, applied at on intermediate code. The main issue raised by this presentation was the applicability of such automatic flow determination to most hard real-time applications. The second presentation dealt with the WCET analysis of object-oriented programs and questioned about the usefulness of the early reclamation of "gain time" (time gained by executing applications quicker than expected) in object-oriented programs. The third presentation introduced a language to express flow information, aimed at being both independent of the methods used to compute WCETs, and independent of the WCET computation tools. It questioned about the possibility to create/agree on a unified language to express flow information, in order to allow the interoperability of different WCET analysis building blocks.

# **Discussions**

The discussions after the three presentations are categorized below according to four broad topics that do not necessarily respect the chronological order in which they appeared.

# On the applicability of automatic flow determination (on today programming languages)

The main interest in the automatic determination of flow information (loop bounds, infeasible paths) is its safety compared to manual annotations, used in most WCET analysis tools. However, automatically determining flow information does not always succeed, since it is equivalent to solving the halting problem. The issue is then to determine whether such automatic methods can be applied to most hard real-time applications.

One of the workshop attendee had experiences with code from avionic and automotive applications showing that the code of such applications is rather simple, and that from 90obtained automatically from the programs source code. Things get more complicated when flow information is obtained from object code, but some information can be made available by the compiler or compiler designer (sizes of arrays, sophisticated translation methods), making the attendee confident on the applicability of automatic flow analysis methods to most hard real-time applications.

One can also distinguish between hard real-time applications, in which the code is intended to be simple, and soft real-time applications, in which complicated flow information occurs more frequently. Automatic flow determination is intended to work better on hard real-time code than on soft real-time code. But this is not a big problem for soft real-time applications, for which the safety of flow information is a less stringent requirement than for hard real-time applications (for the former class of applications, manual annotations or measurements can be used as an alternative to automatic flow determination).

## On the programming discipline required to make flow analysis (and more generally WCET analysis) feasible

Some failures of automatic flow determination techniques can come from code not suited to hard real-time constraints, developed by programmers that are not aware of timing constraints, or programmer that suffer from stringent productivity constraints, like for instance a time-to-market of a few months in the telecommunication area. The point here is that we should not expect the WCET analysis tools to do everything for us. Rather, a programming discipline is required to produce analysable code (programming, especially programming real-time applications, is not easy). This can be seen as a contract between the programmer and the WCET analysis tool: the programmer respects programming discipline such that the WCET analysis tool can obtain WCET automatically. In that respect, WCET analysis tools can help in this educational effort to produce better code, by identifying non real-time code (e.g. non terminating loops).

# On the integration of WCET analysis in the design process of applications

Most WCET analysers assume that the code of the whole application is available, and assume that the application code is correct. This implies that WCET analysis in particular, and that schedulability analysis in general, cannot be used in the very early stages of the design process of applications. This point has been identified as a potential brake on the use of WCET analysis techniques in the industry.

Another issue discussed was the issue of the feedback given by the WCET analysis tools to the system designers. Currently, WCET analysis tools produce a single number (the WCET) as a result of the analysis. Is this the best thing to produce for schedulability analysis? When schedulability analysis fails, is there enough feedback to modify the code or architecture so that the system schedulability can be established? The approach called software performance engineering (Connie Smith) was mentioned; it consists in incrementally specifying and verifying application performance constraints by dividing the application performance constraints into performance budgets that can be specified and verified independently and incrementally.

More generally, the issue of integration of different subproblems related to scheduling has been identified as an important one. Currently, the scheduling problem is traditionally separated into two steps considered independently: schedulability analysis and WCET analysis. But these steps are not actually that separated, especially when hardware components such as caches are used (high worst-case timing cost of an interrupt on an architecture with caches has an impact on the system schedulability). Are we currently using the right design process to separate the scheduling problems into subproblems that can be solved independently?

### On the WCET analysis of object-oriented languages

There seems to be a wish to use object-oriented languages to program real-time applications, especially in the telecommunications area. Some features of object-oriented languages that make their timing analysis difficult or even impossible (or have a too important impact on the performance of applications), have been discussed, and are usually avoided in real-time applications: virtual functions, design patterns, frameworks, dynamic memory management, design factories, ... Two approaches have been discussed to overcome this problem: either using these languages for soft real-time systems, or use them in hard real-time systems and then restrict their capabilities such that their timing behaviour can be analysed, while keeping their expressiveness. In the latter case, estimated WCETs can be more overestimated than in sequential programming languages due to features like dynamic dispatching. In such a situation, identifying and reclaiming gain time as soon as possible can have an important impact on the system overall performance and utilisation.

#### **Other issues**

A recurrent question is why WCET analysis tools are not used in the industry? A potential explanation is the lack of involvement of the compiler vendors in WCET analysis techniques. WCET analysis tools require information on all the things the compiler does (translation schemes, optimisations). Such information is not currently easily available. While there is a need for such information, there is no unified data structure to express such information, and it may be problematic for compiler vendors to provide them since it is a part of their expertise. One can think of uses of WCET analysis with broader applications than hard real-time systems. One can for instance think of using it for comparing the performance of different algorithm implementations.

A last aspect mentioned (not specially linked to high level analysis, but worth mentioning) was the introduction of on-chip multiprocessors. Such architectures blur the traditional separation between WCET analysis and schedulability analysis, because determining the WCET of a task cannot be done independently. It is also the source of observability problems to validate the timing model using in WCET analysis against the actual hardware.

# A Prototype Tool for Flow Analysis of C Programs

Jan Gustafsson, Björn Lisper, Nerina Bermudo, Christer Sandberg, Linus Sjöberg Department of Computer Engineering Mälardalen University, Västerås, Sweden {jgn, blr, nbo, csg}@mdh.se, lsg98020@idt.mdh.se

#### Abstract

We describe a prototype tool for flow analysis. The purpose of the tool is to statically analyse C programs in intermediate code format, and to calculate flow information, like loop bounds. This information will be used by a subsequent low-level analysis to calculate a final worst case execution time. We describe the main steps of the tool, and analyse a simple example to illustrate our method.

# 1 Introduction

Predicting the Worst Case Execution Time (WCET) of programs is an essential step in designing real-time systems, especially hard real-time systems. Methods based on static analysis can guarantee the safeness of the predicted WCET, while measurements, in the general case, can not.

In the presence of loops and recursion, finite iteration bounds must be given to the WCET calculation method. Most often, they are given as *manual annotations* by the programmer. Optional annotations (like information on infeasible paths) may also be given, to reduce the overestimation of the calculated WCET. The annotations can be supplied as comments or in a separate file.

A problem with manual annotations is that the calculation of these are often time-consuming and error-prone. It would be advantageous if these annotations could be calculated automatically. This is the aim of the project described in this paper.

The WCET project is a sub-project within CODER (Cluster on Distributed Embedded Real-Time Systems) in the ASTEC [AST01] competence center. The project consists of two groups, one at Uppsala University (low-level analysis) and one at Mälardalen University in Västerås. The flow analysis research is an activity of the Västerås group.

The flow analysis tool is a part of a planned, complete WCET tool (see  $[EES^+01]$  for details). The flow analysis part will calculate the possible flow of the analysed program. This information will, together with the results from the low-level analysis, be used to calculate a final WCET.

# 2 Overview of the Tool

## 2.1 The Input of the Tool

The tool analyzes C-programs in intermediate code format. We will use the NIC (New Intermediate Code) format (developed within CODER). The full ANSI C language will be supported (including pointers, recursion and unstructured code).

We assume that the code represents a syntactically and logically correct program. For example, we assume that array indices are within bounds. We also assume that the control flow graph of the analysed program is the same as in the final machine code, i.e., that the final steps to machine code does not change the control flow.

Manual annotations are used as a complement when the automatic flow analysis fails.

# 2.2 The Calculation

There will be a possibility to choose between a slower but more exact analysis or a faster but less accurate. It will also be possible to change certain compiler- or system-specific data used in the analysis (like integer type sizes).

## 2.3 The Output

The purpose of the tool is to calculate flow information ("flow facts", see [EE00]), like number of iterations and recursion levels, infeasible paths etc., that will be used in the subsequent low-level analysis. The flow facts are attached to the scope graph [EE00] of the analysed program. A scope graph is a partition of the program into scopes; a scope is a part of the program where certain flow facts are valid.

## 2.4 Flow Analysis Overview

Basically, the analysis of a C program is performed using the steps described in Figure 1.



Figure 1: Basic analysis steps.

- Parser. The C code is parsed to produce a NIC file.
- Pointer analysis. Pointers in the program are analyzed. Information about the resulting points-to sets are stored in the NIC file.
- Optimization. The NIC code is optimized. These first three steps are developed within the WPO project.

- NIC code parser. The optimized NIC code is parsed to produce an internal representation. This internal format is the basis for all subsequent analysis steps.
- An SSA (Static Single Assignment) conversion is performed. The calculated data is added to the existing internal representation.
- Non-conditionals are removed. All assignments to variables that do not affect control flow (transitively) are identified and removed from the program. If all references to a variable are removed, the variable will be removed completely. The reason is to simplify and speed up the rest of the analysis.
- Scope graph construction. The scope graph is constructed using the control flow that can be extracted from the internal representation.
- Syntactical analysis. The code is "scanned" for simple, recognizable loop constructs and the corresponding loop counts are calculated, if possible. The loops are replaced with assignments to the final values for the variables updated in the loop, resulting in a simpler program to analyze in the following step.
- The removal of non-conditionals is run again, since variables may become non-conditional during syntactical analysis.
- Abstract interpretation. The remaining code (after the previous step) is analysed using abstract interpretation. The resulting flow facts are appended to the results file.
- If there are constructs for which the abstract interpretation fails, the user is asked for manual annotations for these. The analysis continues with these two last steps until the complete code is successfully analysed.

# 3 Complete Example

The code below contains a simple and motivating example, activating all the steps of our tool. For simplicity reasons, it does not contain pointers, arrays, or unstructured code. The variable i is assumed to receive a value between 0 and 5 by get\_value().

```
int main(void) {
                                                      int foo(int j) {
 int i, j, k, n = 10, c = 2, p;
                                                        int i, result = 0;
 for (i = get_value(); i <= n; i = i + c) {
                                                       for (i = 0; i < 100; i++) {
   p = i - c * 2; result += 2;
                                                         result += 2;
   if (p) k = i + 2;
                                                         }
                                                       return(result);
   }
 j = i + k;
                                                        }
 j = foo(p);
 return(0);
 }
```

We first parse the code to a NIC file. Conversion to SSA form and removal of non-conditionals  $(j, k, p^1, and result)$  yields a NIC code that is equivalent to the C code in Figure 2.

Next step is to calculate the scope hierarchy below. We see that each function and loop constitutes a scope.



<sup>&</sup>lt;sup>1</sup> The variable p is used to control flow in the if-statement, but it is not used in foo. The analysis therefore removes the later uses of p.

```
int main(void) {
    int i, n = 10, c = 2, p;
    int i;
    for (i = get_value(); i <= n; i = i + c) {
        p = i - c * 2;
        if (p) {}
    }
    }
    foo(0);
    return(0);
}</pre>
```

Figure 2: Example program after removal of non-conditionals

The syntactical analysis will recognize the loop in foo as analyzable and output the flow fact

foo\_L1:[]: $x_{header(foo_L1)} = 100$ 

which means that the loop in scope **foo\_L1** iterates exactly 100 times. The notion  $x_{header(foo_L1)}$  refers to the iteration count of the loop header. The function **foo** will be changed by the syntactical analysis as shown below. We see that the loop has been replaced by an assignment.

```
int foo(int j) {
    int i = 100;
    return(0);
}
```

A new run of removal of non-conditionals removes the variable i in foo since it does not affect the control flow.

Abstract interpretation of the remaining program will yield the following flow facts for the remaining loop:

- 1. main\_L1:[]: $x_{\text{header}(\text{main}_{L1})} \geq 3$
- 2. main\_L1:[]: $x_{\text{header}(\text{main}_{L1})} \leq 6$
- 3. main\_L1:<4..6>: $x_{true} = 1$

The first two flow facts means that the loop in main iterates between 3 and 6 times. The second means that the program will always take the true edge in the if-statement in iterations 4 to 6 of the loop.

# References

- [AST01] ASTEC (Advanced Software TEChnology) WWW Homepage. URL: http://www.astec.uu.se/, November 2001.
- [EE00] J. Engblom and A. Ermedahl. Modeling complex flows for worst-case execution time analysis. In Proc. 21<sup>st</sup> IEEE Real-Time Systems Symposium (RTSS'00), November 2000.
- [EES<sup>+</sup>01] J. Engblom, A. Ermedahl, M. Sjödin, J. Gustafsson, and H. Hansson. Worst-case execution-time analysis for embedded real-time systems. Springer International Journal of Software Tools for Technology Transfer, (STTT), 2001.
- [EG97] A. Ermedahl and J. Gustafsson. Deriving Annotations for Tight Calculation of Execution Time. In Proc. 3<sup>rd</sup> International European Conference on Parallel Processing, (Euro-Par'97), LNCS 1300, pages 1298-1307, August 1997.

# A Novel Gain Time Reclaiming Framework Integrating WCET Analysis for Object-Oriented Real-Time Systems

Erik Yu-Shing Hu<sup>\*</sup>, Andy Wellings and Guillem Bernat

Real-Time Systems Research Group Department of Computer Science University of York, York, YO105DD, UK E-mail: {erik,andy,bernat}@cs.york.ac.uk

## Abstract

This paper proposes a novel gain time reclaiming framework integrating WCET analysis for objectoriented real-time systems in order to provide greater flexibility and without loss of the predictability and efficiency of the whole system. In this paper we present an approach which demonstrates how to improve the utilisation and overall performance of the whole system by reclaiming gain time at run-time. Our approach shows that integrating WCET with gain time reclaiming not only can provide a more flexible environment to develop object-oriented real-time applications, but it also does not necessarily result in unsafe or unpredictable timing analysis.

**Keywords :** Gain Time, Real-Time Java, Worst-Case Execution Time (WCET) Analysis, Object-Oriented WCET

### 1 Introduction

There is a trend towards using object-oriented programming languages, such as Java and C++, to develop real-time applications. The success of hard realtime systems, undoubtedly, relies upon their capability of producing functionally correct results within defined timing constraints. In order to achieve this, the processor and resource requirements of the hard real-time tasks have to be reserved. However, this may result in under utilisation and lead to very poor performance for aperiodic tasks. Unfortunately, object-oriented programming languages support more dynamic behaviour than procedural programming languages, and some of these features may bring about object-oriented applications having a more pessimistic worst-case behaviour. In consequence, object-oriented real-time systems may suffer from significantly lower utilisation and poorer overall performance of the whole system than procedural realtime systems.

Most scheduling algorithms assume that the WCET estimation of each task is known prior to doing the schedulability analysis. Typically, the WCET analysis and schedulability analysis are carried out sepa-Sophisticated techniques [6, 16, 17, 19], are rately. used in WCET analysis, for instance to model caching and pipelining, to achieve safe and tight WCET estimation. However, most WCET analysis approaches are only considered in relation to procedural programming languages. Performing WCET analysis on objectoriented programs must take into account additional dynamic features, such as dynamic dispatching and memory management. Some research groups have proposed various approaches [10, 18] to address these issues, but most approaches result in developing environments which are inflexible and very limited.

In contrast with the WCET analysis, a number of research groups have proposed various flexible scheduling algorithms [5, 14], for instance priority server algorithms [5] and slack stealing algorithm [14], to provide a more flexible real-time development environment with greater performance of the whole system. In general, these flexible scheduling algorithms are mainly focused on the use of WCET estimation to improve the performance of the aperiodic tasks at run-time. They have, however, paid insufficient attention the fact that, for the most part, hard real-time tasks are not executing via the worstcase execution time path. Therefore, even though they have demonstrated very complex scheduling algorithms to improve the average performance of the whole system, the improvements are still limited and the overhead

<sup>\*</sup>This work has been funded by the EPSRC under award number GR/M94113.

of the implementation is extremely high or it is sometimes not even possible to implement them in practice.

On the whole, the spare capacity of the real-time system may be divided into three groups [8]: extra capacity, gain time, and spare time. Extra capacity is the capacity which is not allocated for hard real-time tasks during the design phase. This can be identified off-line. The gain time is produced when the hard real-time tasks execute in less than their worst-case execution time estimations. This may only be reclaimed at run-time since it depends on the actual executions of task [8]. The spare time may be defined as a situation in which the sporadic tasks do not arrive at their maximum rate. Most flexible scheduling algorithms are mainly focused on reclaiming the extra capacity of the system. Only a few research approaches [11, 9, 1] have discussed how to reclaim the gain time. However, they have tended to focus on procedural programming languages, rather than on object-oriented programming languages.

We have proposed a static timing analysis environment for the Java programming language [13] and an approach [12] to address dynamic dispatching issues for object-oriented real-time systems. However, these only take into account hard real-time threads. To provide a more flexible computation model, one needs to consider how to improve the performance of aperiodic tasks without rendering hard real-time tasks unsafe.

In this paper we propose an approach which demonstrates how to improve the utilisation and overall performance of the whole system by reclaiming gain time at run-time. The primary focus of this paper is to demonstrate how to reclaim the gain times from periodic realtime tasks, rather than how to apply them in scheduling algorithms. Integrating the gain time reclaiming with the scheduling algorithm is outside the scope of this paper. Techniques such as Dual-Priority Scheduling [7] and Dynamic Sporadic Server [5] are applicable. We use a gain time reclaiming mechanism to compensate for the tradeoff among flexibility, efficiency and predictability. In our approach, the predictability of hard real-time tasks is strengthened during the design phase and the performance of the whole system is reinforced with gain time reclaiming during run-time. It shows that integrating WCET analysis with gain time reclaiming not only may achieve high utilisation and high performance of the whole real-time system, but also keep the flexibility and reusability of the object-oriented real-time applications. The major contributions of this paper are:

- presenting how to address the dynamic behaviour of object-oriented programming features with minimum annotations
- demonstrating how to reclaim the gain time of high

performance object-oriented real-time systems with the gain time reclaiming graphs

• balancing the flexibility and predicability of objectoriented real-time applications by integrating WCET analysis

The rest of the paper is organised as follows. A survey of related work and an overview of our previous work are given in Section 2 and Section 3 respectively. Section 4 demonstrates how gain time can be reclaimed in objectoriented real-time systems. Finally, the conclusion and future work are presented in Section 5.

## 2 Related Work

This section gives a survey of the related work of gain time analysis [11, 9, 1] in the literature. Haban and Shin have proposed an approach [11] placing software triggers at the end of basic blocks in task code to measure actual execution time. In [11], comparing the actual execution time which calculated at the software triggers point with pre-determined WCET values, the gain time of the specific basic block can be calculated. In a similar way, Dix et al. have proposed an approach [9] adding *milestones* into task code to calculate the maximum remaining execution time of the particular task. However, both approaches reclaim the gain times after they have been generated and do not integrate with WCET analysis.

Audsley et al. [1] have introduced a gain point mechanism to reclaim gain times of the basic blocks of a task code as early as possible. In [1], the use of gain point can be grouped into four separate forms, including static gain point for static code, dynamic gain point for loop constructs, efficiency gain point for detecting hardware speed-ups, and resource usage gain point for identifying spare resources. Yet, Audsley et al.'s approach and the previous two approaches do not take into account object-oriented programming features and gain times resulting from functional constraints impacting on the program's execution.

## 3 Previous Work

Our previous work, called Extended Real-Time Java [13] (XRTJ), extends the current Real-Time Java architecture [4] proposed by the Real-Time Java Expert Group. The XRTJ architecture has been developed with the whole software development process in mind: from the design phase to run-time phases. For example, using our approach, the system can be evaluated during the design, and the timing constraints of the application can be validated during run-time. We integrate our approach with portable WCET analysis, proposed by Bernat et al. [3] and extended by Bate et al. [2], for the WCET estimation. The portable WCET analysis uses a three-step approach: high-level analysis (i.e. analysing the Java programs), low-level analysis (i.e. performing platform-dependent analysis on Java byte code instructions implemented for the target platform), and conducting the combination of the high-level analysis with the low-level analysis to compute the actual WCET bound of the analysed code sections.

In our previous approach [13], we have introduced the *Extensible Annotations Class* (XAC) format, which stores extra information that cannot be expressed in the source code. The XAC format is an annotation structure that can be stored in files or as an additional code attribute in *Java Class Files* (JCF). We have also addressed dynamic dispatching issues in object-oriented real-time applications [12]. Here, minimum annotations are provided to ensure the predictability of dynamic binding methods and to estimate safe and tight WCET for hard real-time applications. However, our previous work mainly focused on analysing hard real-time objectoriented tasks.

#### 3.1 WCET Annotations

This section reviews three annotations proposed in our previous work [12] to apply them in the rest of this paper. We have introduced //@UseWCET() to address dynamic dispatching problems of hard real-time tasks in object-oriented real-time systems. This annotation is applied to denote a specific method in the applications and specifies the worst case value for execution of the method. We have also introduced //@DefineScope() annotation to provide for defining a simple or nested scope of program. This annotation may enable the WCET analysis to indicate the complicated structure of the applications.



#### Figure 1. An example of //@maxWCET() annotation

Furthermore, to address more complicated class hierarchies where there is more than one particular method that may be invoked, the //@maxWCET() annotation is introduced. It can suggest that the WCET of a dispatching method should be considered to be the maximum WCET of the class family<sup>1</sup> containing that method. Subsets of the class family can also be specified. In this annotation, "&" may be applied to denote the whole class family of a class. In addition, "+" and "-" can also be used to express the union or subtraction of a single class or a class family for the method. An example of the use of //@maxWCET() annotations is given in Figure 1.

### 4 Gain Time Reclaiming

In order to balance the tradeoff among the flexibility, predictability and efficiency of the real-time systems, gain time reclaiming needs to be applied. For the most part, the gain time reclaiming in object-oriented programming languages may be classified in three sections: structural constraints reclaiming, functional constraints reclaiming, and object constraints reclaiming. Further details of each reclaiming mechanism are discussed below. Note that the WCET annotations used in the following examples to discuss the gain time reclaiming mechanism can be added either manually by developers or automatically by modified compilers or tools.

#### 4.1 Structural Constraints Reclaiming

From the point of view of the syntax of the programming languages, the real-time tasks allow construction with a number of basic blocks, conditional branches, and call procedures. These components of a real-time task, in general, may be represented by a *control flow* graph (CFG). It can be observed that the actual execution time of real-time tasks may vary, if the execution paths of the task or iteration times are varied at runtime. This section is mainly concerned with reclaiming gain times, which depend on the structural constraints of a specific real-time task.

Our approach is similar to Audsley et al.'s approach [1], which is proposed for procedural programming language. We have defined two annotations (A1 & A2), which are given in Table 1, to cope with structural constraints reclaiming. On the whole, the static gain time, such as pre-calculated units or paths, can be annotated with annotation A1, and the dynamic gain time, defined for unknown iteration times, can be interpreted with annotation A2. Essentially, pre-calculated units can be represented with machine cycles of the target

 $<sup>^1{\</sup>rm A}$  class family of a class is a set of the classes including the class itself and all the child classes inherited from it.

//@ GainTime(Units /Vectors /path /mode /method)	 A1
$//@$ dynGainTime(maxLoopcount, Scope_Name)	 A2
$//@$ objGainTime( $Object\_Name$ )	 A3

#### **Table 1. Gain Time Reclaiming Annotations**

machine if the source code of the application is translated into machine code directly. However, it could be difficult to estimate the exact machine cycles of Java applications because of the portability of Java architecture. In this case, the concept of Worst-Case Execution Frequency (WCEF) vector [2], which represents execution-frequency information about basic blocks and more complex code structures that have been collapsed during the high level timing analysis phase, may be used instead of pre-calculated units. These WCEF vectors may be used to calculate the exact gain time when the information about the target machine is available. In order more easily to understand our approach with a straightforward example, the machine cycle units are used in the rest of paper.

```
1
^{2}
   public check_data() {
     int i, morecheck, wrongone;
 3
     i=0; morecheck=1; wrongone=-1;
 4
     //@ DefineScope(checkLoop)
 6
     while (morecheck) {
 8
9
       if(data[i] < 0) \{ //Say WCET=20 cycles
10
         //@ GainTime(100); (i.e. 120-20 cycles)
         wrongone=i; morecheck=0;
11
         //@ GainTime(Error_mode);
12
13
       else { //Say WCET=120 cycles
14
15
         if(++i \ge DATASIZE)
16
           morecheck=0;
17
18
19
     //@ dynGainTime(50, checkLoop);
^{20}
^{21}
22
     if (wrongone >= 0) { //Say WCET=10 cycles
^{23}
       //@ Mode(Error_mode);
^{24}
^{25}
       return 0;
^{26}
27
             //Say WCET=50 cycles
     else
^{28}
       //@ Mode(Noml_mode);
^{29}
30
       return 1;
31
32
     }
33 }
34
```

# Figure 2. An example of gain time reclaiming [15]

As shown in Table 1, annotation A1 may provide

various types of parameter, such as units (i.e. machine cycles), vectors (i.e. WCEF vectors), paths, modes and methods. As shown in Figure 2, the if-then-else basic block can reclaim 100 cycles at Line 10, if the condition expression is TRUE (i.e. data[i]<0) and the while-loop is part of its worst case path. With respect to the dynamic gain time, we can simply add an annotation to a non-constant iteration loop, such as for-loop or while-loop, in order to reclaim gain times at runtime.

Essentially, the gain time can be reclaimed as soon as the exact execution path of the task or iteration time are identified. One should note that the dynamic gain time reclaiming needs to be provided with the maximum loop bound. It may also note that either the runtime system, such as the Virtual Machine, must support a mechanism to count the exact iteration of the loop at run-time or additional code must be introduced by an annotation aware compiler to count the loops. The structural constraints reclaiming of a specific task may be represented with a Structural Gain Time Reclaiming Graph (SGTRG), which illustrates the exact places (i.e. offset number of the machine code or Java byte code) and amounts (i.e. machine cycles or WCEF vectors) of gain time that may be reclaimed. Note that it could be possible that the actual reclaimed gain time is less than the run-time overhead of the reclaiming. In this situation, the gain time should be either neglected or accumulated until it is worth reporting.

#### 4.2 Functional Constraints Reclaiming

This section is mainly concerned with reclaiming the gain times which suffer from functional constraints. This covers the issues that remain from the previous sections which did not take into account the functional and data dependencies of the exclusive paths or modes of the real-time task.

Identifying the exclusive paths [15] or various modes [6] in order to calculate the WCET estimation of the real-time program is widely used in the WCET field. Based on design knowledge, the annotations of the exclusive paths or modes may be distinguished during the design phase. Using these annotations, the WCET estimation of each exclusive path or mode may be calculated. However, one should note that it is possible that the WCET estimations of the exclusive paths or different modes are spread over a wide range, and the exact execution path or mode cannot be determined during the design phase. As a result, the WCET estimation could be very pessimistic. In order to address this, we propose a gain time reclaiming framework which takes into account the functional constraints of the programs.

In our approach, we use the gain time annotation (A1), given in Table 1, to identify where the exclusive path or mode can be determined. As soon as the specific execution path or mode is determined or executed, the associated gain time of the executed path or mode can be reclaimed. Again using the previous example in Figure 2, the A1 annotation can be annotated at Line 12 to reclaim the functional associated gain time at run-time. It can be observed that using functional constraints reclaiming may reclaim the gain time earlier than the structural constraint reclaiming. The functional constraints reclaiming of a specific task may be represented with a *Functional Gain Time Reclaiming Graph* (FGTRG), which illustrates the exact places and amounts of gain time that may be reclaimed.

#### 4.3 Object Constraints Reclaiming

So far, we have only discussed the gain time reclaiming which may apply to both procedural and objectoriented programming languages. We have argued for the need to use dynamic dispatching and demonstrated how to guarantee the deadline of hard real-time tasks in our previous work [12]. Our previous approach has shown that allowing the use of dynamic dispatching not only can provide a more flexible way to develop objectoriented hard real-time applications, but it also does not necessarily result in unpredictable timing analysis. Essentially, a //@maxWCET() annotation is used to indicate the WCET of a dynamic dispatching method call. However, we cannot avoid the fact that the use of //@maxWCET() might have relatively pessimistic results if the class family is too large or the WCET estimations for different classes are spread over a wide range. In order to compensate for the penalty of the flexibility of the object-oriented programming, gain time reclaiming is required.

Before discussing further details of the object constraints reclaiming, two technical terms are introduced below.

• An Object Type Lifetime Graph (OTLG) is a diagram which represents lifetimes of types of particular objects in a specific task. An OTLG is made of two types of component: node and edge. A node denotes a place where the type of the object is changed, whereas an *edge* illustrates the lifetime of a particular type of object between two nodes.

• An Object Gain Time Reclaiming Graph (OGTRG) is a diagram which illustrates places where the object constraint reclaiming may take place. The OGTRG can be produced from the analysis of the OTLG and CFG of a specific task.

```
Assume that Class A is a parent
class. Class B, C and D extend A,
and override the m1() methd.
*****
class App extends RealtimeThread {
 public void run() {
   //@ objGainTime(aa);
   A aa= \mathbf{new} A();
   //@ objGainTime(bb);
   B bb = new B();
   C cc = new C();
   D dd = new D();
   Initial values of x, y and z
   are from the environment.
   if(x > 5) {
    cc = dd;
   }
   if(y == 5) \{
    aa = dd;
   else {
    aa = bb;
   }
   // type changing
   bb = cc;
   if(z == true) {
    aa.m1;
    aa m1.
   } else {
    aa.m1;
   }
   bb.m1;
   bb.m1;
 }
```

#### Figure 3. An example of object gain time reclaiming

Essentially, the value of the dynamic dispatching gain time of each object can be calculated as follows: //@GainTime() = //@maxWCET()-//@UseWCET(). The



Figure 4. A diagram of producing OGTRG

annotations of the object gain time reclaiming may be generated by using design knowledge or by producing an OGTRG. In order to reduce the run-time overhead, annotation A3 may be applied to define which object's gain times are going to be reclaimed. Considering the example in Figure 3, two //@objGaintime() annotations are annotated in the run() method. This means that only these two objects will be taken into account while producing the OGTRG of the run() method. The procedure of object gain time reclaiming is given as follows.

The CFG can be produced from the source code (or Java class file) for each hard real-time task. Based on the CFG, an OTLG for each object or those objects denoted with annotation A3 in the real-time task can be produced. In the OTLG, symbolic references may be applied to represent the relationship between the dynamic dispatching objects of the same class family during runtime. Using the CFG and the OTLG of each object, the exact places and amounts of gain time reclaiming can be identified. These gain time reclaiming places can be illustrated in the OGTRG for each object. Following this, the gain time reclaiming of all objects in the real-time task can be merged together and provided for the runtime environment (or Java virtual machine) to reclaim them. A diagram which illustrates the transformation from CFG to OGTRG is given in Figure 4.

Solving the symbolic expression of an associated class family can improve the reclaiming as early as possible. As shown in figures 3 and 4, the gain time of the object bb can be reclaimed as soon as the type of the object cc is determined.

### 5 Conclusion and Future Work

This paper has demonstrated a novel gain time reclaiming framework integrating WCET analysis for object-oriented real-time systems. Our approach shows that integrating WCET with gain time reclaiming not only can provide a more flexible environment to develop object-oriented real-time applications, but may achieve high utilisation and high performance of the whole realtime system.

Here, we have mainly discussed the dynamic behaviour of object-oriented features which is exclusively restricted to a consideration of the language syntax and semantic aspects. In order to cover as much dynamic behaviour of the object-oriented programming features as possible, our future work has to take into account: memory management, dynamic loading and extension, and remote method invocation (RMI) issues.

#### References

- N. C. Audsley, R. I. Davis, and A. Burns. Mechanisms for Enhancing the Flexibility and Utility of Hard Real-Time Systems. In Proc. of the 15th IEEE Real-Time Systems symposium (RTSS), pages 12–21, December 1994.
- [2] I. Bate, G. Bernat, G. Murphy, and P. Puschner. Low-Level Analysis of a Portable Java Byte Code WCET Analysis Framework. In 6th IEEE Real-Time Computing Systems and Applications (RTCSA2000), pages 39–48, December 2000.
- [3] G. Bernat, A. Burns, and A. Wellings. Portable Worst-Case Execution Time Analysis Using Java Byte Code. In proc. 6th Euromicro conference on Real-Time Systems, pages 81–88, June 2000.
- [4] G. Bollella, J. Gosling, B. M. Brosgol, P. Dibble, S. Furr, D. Hardin, and M. Turnbull. *Real-Time Specification for Java*. Addison Wesley, 2000.
- [5] G. Buttazzo. Hard Real-Time Computing Systems: Predictable scheduling algorithms and applications. Kluwer Academic Publishers, 1997.
- [6] R. Chapman, A. Burns, and A. Wellings. Integrated Program Proof and Worst-Case Timing Analysis of SPARK Ada. In Proc. of the Workshop on Language, Compiler, and Tool Support for Real-Time Systems, June 1994.
- [7] R. Davis and A. Wellings. Dual Priority Scheduling. In Proceedings of 16th IEEE of Real-Time Systems Symposium (RTSS), pages 100–109, December 1995.
- [8] R. I. Davis. On Exploiting Spare Capacity in Hard Real-Time Systems. Ph.d. thesis, Department of Computer Science, University of York, UK, July 1995.
- [9] A. Dix, R. Stone, and H. Zedan. Design Issues for Reliable Time-Critical Systems. Technical Report YCS-133, Department of Computer Science, University of York, UK, 1990.
- [10] J. Gustafsson. Analysing Execution Time of Object-Oriented Programs with Abstract Interpretations. Ph.d. thesis, Department of Computer Systems, Information Technology, Uppsala University, Sweden, May 2000.
- [11] D. Haban and K. Shin. Application of Real-Time Monitoring to Scheduling Tasks with Random Execution Times. *IEEE Transactions on Software Engineering*, 16(12), December 1990.
- [12] E. Y.-S. Hu, G. Bernat, and A. J. Wellings. Addressing Dynamic Dispatching Issues in WCET Analysis for Object-Oriented Hard Real-Time Systems. Proceedings of the 5th IEEE International Symposium on Object-Oriented Real-Time Distributed Computing ISORC-2002, pages 109–116, April 2002.
- [13] E. Y.-S. Hu, G. Bernat, and A. J. Wellings. A Static Timing Analysis Environment Using Java Architecture for Safety Critical Real-Time Systems. Proceedings of the 7th IEEE International Workshop on Object-Oriented Real-Time Dependable Systems WORDS-2002, pages 77–84, January 2002.

- [14] J. Lehoczky and S. Ramos-Thuel. An Optimal Algorithm for Scheduling Soft-Aperiodic Tasks in Fixed-Priority Preemptive Systems. In Proceedings of 13th IEEE of Real-Time Systems Symposium (RTSS), pages 110–123, December 1992.
- [15] Y. Li and S. Malik. Performance Analysis of Embedded Software Using Implicit Path Enumeration. ACM SIG-PLAN Workshop on Language, Compilers and Tools for Real-Time Systems, June 1995.
- [16] S. Lim, Y. Bae, G. Jang, B. Rhee, S. Min, C. Park, H. Shin, K. Park, and C. Kim. An accurate worst case timing analysis for RISC processors. *IEEE Transactions on Software Engineering*, 21(7):593–604, July 1995.
- [17] F. Mueller. Static Cache Simulation and its Applications. Ph.d thesis, Department of Computer Science, Florida State University, July 1994.
- [18] P. Persson and G. Hedin. An Interactive Environment for Real-Time Software Development. Proceedings of the 33rd International Conference on Technology of Object-Oriented Languages (TOOLS Europe 2000), June 2000. St. Malo, France.
- [19] P. Puschner and A. Burns. A Review of Worst-Case Execution-Time Analysis. *Real-Time Systems*, 18(2/3):115–128, 2000.

# A Unified Flow Information Language for WCET Analysis

Andreas Ermedahl<sup>†</sup> IT-Dept. Uppsala University Box 337, SE-751 05 Uppsala Sweden andreas.ermedahl@it.uu.se Jakob Engblom<sup>†</sup> IAR Systems AB Box 23051, SE-750 23 Uppsala Sweden jakob.engblom@iar.se Friedhelm Stappert\* C-LAB Fürstenallee 11, 33102 Paderborn Germany friedhelm.stappert@c-lab.de

#### Abstract

In this paper we raise the question if it is possible to create a unified flow information language that all WCET research groups can agree upon, and that is independent of flow analysis and calculation methods.

We discuss desired characteristics of such a flow information language and describe the type of flows that it should be able to express. We present our previously published flow fact annotation language and discuss how it fulfils the desired language properties.

#### 1. Introduction

A correct WCET calculation method must take into account the possible program flow, like loop iterations and function calls. For expressing program flows numerous annotation languages have been presented in the WCET literature. The expressiveness and the type of flows that can be handled by these languages mostly depend on the characteristics of flow analysis methods used, rather than being targeted for the potential WCET tool user.

To generate a WCET estimate, we consider a program to be processed through the phases of *program flow analysis*, *low level analysis* and *calculation*. Most WCET research groups make a similar division notationally, but sometimes integrate two or more of the phases into a single algorithm.

The program flow analysis phase determines possible program flows, and provides information about which functions get called, how many times loops iterate, if there are dependencies between if-statements, etc. The information can be obtained by manual annotations (integrated in the programming language [14] or provided separately [6, 9, 19]). The flow information can also be derived using automatic flow analysis methods [7, 10, 13, 22].

In the calculation phase a program WCET estimate is derived, combining the information derived in the program flow and low-level analysis phases. There are three main categories of calculation methods proposed in literature: *tree-based*, *path-based*, and *IPET* (Implicit Path Enumeration Technique).

In a *tree-based* approach the WCET is calculated in a bottom-up traversal of a tree generally corresponding to a syntactical parse tree of the program, using rules defined for each type of compound program statement (like a loop or an *if*-statement) to determine the execution time at each level of the tree [1, 2, 16, 20].

In a *path-based* approach the possible execution paths of a program or piece of a program are explored explicitly to find the longest path [10, 12, 22, 23]. The path-based approach is natural within a single loop iteration or function.

In *IPET*, program flow and low-level execution time are modeled using arithmetic constraints [6, 9, 15, 18, 21]. Each basic block and program flow edge in the program is given a time  $(t_{entity})$  and a count variable  $(x_{entity})$ , and the goal is to maximize the sum  $\sum_{i \in entities} x_i * t_i$ , subject to constraints reflecting the structure of the program and possible flows.

#### 2. Representing Program Flow

The program flow phase can be further divided into three different subphases:

- 1. Flow analysis: Obtaining flow information. By manual annotations or automatic flow analysis.
- 2. Flow representation: Representing the results of the flow analysis.

<sup>&</sup>lt;sup>†</sup> This work is performed within the Advanced Software Technology (ASTEC, http://www.docs.uu.se/astec) competence center, supported by the Swedish National Board for Industrial and Technical Development (NUTEK, http://www.nutek.se).

<sup>\*</sup> Friedhelm is a PhD student at C-LAB (www.c-lab.de), which is a cooperation of Paderborn University and Siemens.

3. Calculation: Using the control flow information (as represented in the flow representation) in the final WCET calculation.

Some WCET methods integrate two or more of the phases. We believe that the separation of the flow analysis from the calculation reduces the complexity of each stage. Also, by keeping the flow analysis phase separate from the flow representation, results from several different flow analysis methods and manual annotations can be integrated and used together in the calculation phase.

When designing a language for expressing flow information there are a number of choices to be made:

- *Expressiveness:* What type of flows should be possible to express? What type of language constructs should be used?
- *Code relation:* How is the information related to different entities in the program code?
- *Calculation conversion:* How should the information be used in the final calculation phase?

#### 2.1. Expressiveness

We first note, that a natural way to give flow information is by constraining the number of times different program entities, e.g. loops, statement, nodes or edges, can be taken. This can either be precise bounds, e.g. that a loop is iterated exactly ten times, or upper or lower bounds, e.g. that node A can't be taken more than five times. It is also beneficial if we can relate the executions of different program entities, e.g. that node A and node B will always be executed together.

The language can consist of named special relations between entities (e.g. using constructs like Parks samepath(A,B) and nopath(A,B) [19]). An alternative is to use a more generic style based on math, like our flow fact language [6]. The benefit of a generic math-based language is that it can express flows that are hard to put in words and that there is no obvious limit to the types of flows that can be expressed. On the other hand, a special purpose language is easier to understand, but requires that new language constructs are invented in order to express new flows.

The language must reflect the flows found in realworld programs. Researchers have investigated embedded software [4], the RTEMS operating system [3] and common signal-processing algorithms [8]. The results are not in complete agreement on the properties and flows typical for embedded software, showing that more research and knowledge is needed here.

One observation is that flow information is mostly local in its nature, specifying something valid for a small part of a program or a particular invocation of a function. Thus, it is not always suitable to specify flow information once for each entity in the program. E.g. we would like to be able to specify that some node A can't be executed during the first five iterations of a loop or give a loop bound valid for just some particular executions of a loop. A language should allow for such local flow information to be expressed.

#### 2.2. Code Relation

First we note that it is natural to express flow information in relation to the entities available in the program code. Flow information can be provided in relation to the source code, intermediate code in a compiler, or the object code. If provided on source code level, the information must be mapped to the object code to be used in the WCET calculation. In the presence of optimizing compilers, this problem is nontrivial [5, 17].

Automatic flow analysis is probably easier to perform at the source code or intermediate code, since variables and other entities of interest are harder to identify in optimized object code. Also, for the potential WCET-tool end-user manual annotations are typically easier to provide at the source-code level.

Another issue is if the flow information should be included as a part of the programming language or provided outside the program. The benefit of language inclusion is that it forces the programmer to write code in an analysable manner. However, this requires compiler support and makes it harder to try different scenarios.

Specifying the flow information outside the program source allows it to free itself from the static structure of the program. For example, by using a *call-graph* representation, we can differ between invocations of the same function when called from different places in the code. An example of the extended version is our *scope graph* represention [6].

A good language should provide *stability* in that program changes not related to annotated code should not force the annotations to change. For example, a problem with expressing flow information on the object code level is that the information might need to be regenerated every time the program code changes.

An important issue is the ability to handle *unstructured code*, e.g. due to uses of **goto** and jumps into loops. An optimizing compiler might produce unstructured object code from structured source code, and automatic code for state machines also tends to be unstructured. A general purpose flow information language must be general enough to express flows over such unstructured code.

#### 2.3. Calculation Conversion

Regardless of the flow information language used the extracted flow information must be "compiled" or



Figure 1. Example of Code with Different Type of Flows

"adapted" to the calculation method used. The adaptation must be *safe*: never exclude execution paths which are considered possible by the flow information, and *tight*: including as few extra execution paths compared to the provided flow information. Figure 1 gives example code showing that not all calculation methods can take advantage of all types of flow information.

The tree-based method [1, 2, 16, 20] is conceptually simple and computationally cheap, but has problems handling flow information, since the computations are local within a single program statement and thus cannot consider dependencies between statements. For example, the code and flow information in Figure 1(a) causes problems in a tree-based calculation method since the timing of the first if-statement will be calculated in isolation from the second if-statement.

The path-based approach is natural within a single loop iteration or other executions of one loop [11, 23]. The method has problems with flow information stretching over loop borders and/or flow information on the *total* number of times entities are taken. For example, the path-based method has problems handling the "triangular" loop dependency in Figure 1(b). If WCET calculation is performed locally, the WCET calculation for the inner loop will assume 10 iterations, and the WCET calculation for the outer loop will use 10 executions of the inner loop, leading to the body of the inner loop being counted 100 times, when it is actually never executed more than 55 times.

For IPET very complex flows can be expressed using constraints, but all flow information needs to be given on a global program level [6, 9, 15, 18, 21]. This contradicts the need to specify flow information in a local context. As shown in [6], local flows can be handled by unrolling the program and lifting the information to a global level. Since flow information is given as relations over count variables some type of flow implications are problematic to express. E.g. Figure 1(c) shows an example of code where we would like to express an implication dependency like: "if F is taken once then (and only then) G can be taken several times, but if F is not taken then G can not be taken either".

#### 3. Our Flow Fact Language

This chapter describes our previously published flow fact annotation language [6] and discusses how it fulfils the desired language properties.



Figure 2. Scopes with Attached Flow Facts

The program representation used is the *scope graph*. It is a hierarchical representation of the dynamic structure of the program. Each *scope* corresponds to a certain repeating or differentiating execution context in the program, e.g. loops and function calls, and describes the execution of the object code of the program within that context. Figure 2(b) shows the scope graph generated for the code in Figure 2(a).

A scope consists of a number of *nodes* and *edges*. A node belongs to exactly one scope, and represents the execution of a certain basic block in the program in the environment given by the scope and its ancestors. For each scope, a header node must be given. If the scope iterates, each iteration must pass the header node, and a bound on the number of iterations has to be provided.

To express more complex program flow information than just basic loop bounds each scope can carry a set of *flow facts* [6]. The flow facts use constraints local to a scope to describe the flow. The constraints can be given for a range of iterations, or all iterations of a certain loop. They can also be local within a single iteration ("foreach facts") or represent a total over all iterations ("total facts").

The scope graph in Figure 2(b) has been decorated with some flow facts.

Flow fact inner:<>: $x_{\rm C} + x_{\rm F} \leq 1$  is a foreach fact and gives that the nodes C and F cannot be executed on the same iteration of the scope inner (an infeasible path), while the flow fact inner:<6..10>: $x_{\rm C} = 0$  gives that for each entry of inner, during iterations 6 to 10 of inner, node C can not be executed.

Flow fact inner: [1..10]:  $x_{G} = 3$  is a total fact that gives that, for each entry of inner, during the ten first

iterations, node G must be taken exactly three times.

Compared to the criteria given above, we note that the flow facts language uses the math-based style and allows us to give local information. The information is given outside the code and uses an expanded version of the call graph (and thus the control flow graph). In its current version, it cannot handle all types of unstructured code due to the need for a header, and since it relates to the object code, it is very sensitive to program changes.

It has been used to perform both IPET- and pathbased calculations [6, 23], but not all facts could be used in the path-based approach. It is interesting that the path-based calculation recognized certain types of facts as meaning "samepath" or "not samepath", and exploited these by rewriting the graph.

#### References

- R. Chapman. Program Timing Analysis. Dependable Computing System Centre, University of York, England, May 1994.
- [2] A. Colin and I. Puaut. Worst Case Execution Time Analysis for a Processor with Branch Prediction. *Journal of Real-Time Systems*, May 2000.
- [3] A. Colin and I. Puaut. Worst-Case Execution Time Analysis for the RTEMS Real-Time Operating System. In Proc. 13<sup>th</sup> Euromicro Conference of Real-Time Systems, (ECRTS'01), June 2001.
- [4] J. Engblom. Static Properties of Embedded Real-Time Programs, and Their Implications for Worst-Case Execution Time Analysis. In Proc. 5<sup>th</sup> IEEE Real-Time Technology and Applications Symposium (RTAS'99). IEEE Computer Society Press, June 1999.
- [5] J. Engblom, P. Altenbernd, and A. Ermedahl. Facilitating worst-case execution times analysis for optimized code. In *Proc. of the 10<sup>th</sup> Euromicro Workshop of Real-Time Sys*tems, pages 146–153, June 1998.
- [6] J. Engblom and A. Ermedahl. Modeling Complex Flows for Worst-Case Execution Time Analysis. In Proc. 21<sup>th</sup> IEEE Real-Time Systems Symposium (RTSS'00), November 2000.
- [7] A. Ermedahl and J. Gustafsson. Deriving Annotations for Tight Calculation of Execution Time. In Proc. Euro-Par'97 Parallel Processing, LNCS 1300, pages 1298–1307. Springer Verlag, August 1997.
- [8] R. Ernst and W. Ye. Embedded program timing analysis based on path clustering and architecture classification. In International Conference on Computer-Aided Design (IC-CAD '97), 1997.
- [9] C. Ferdinand, F. Martin, and R. Wilhelm. Applying Compiler Techniques to Cache Behavior Prediction. In Proc. ACM SIGPLAN Workshop on Languages, Compilers and Tools for Real-Time Systems (LCT-RTS'97), 1997.
- [10] C. Healy, R. Arnold, F. Müller, D. Whalley, and M. Harmon. Bounding Pipeline and Instruction Cache Performance. *IEEE Transactions on Computers*, 48(1), January 1999.

- [11] C. Healy, M. Sjödin, V. Rustagi, D. Whalley, and R. van Engelen. Supporting timing analysis by automatic bounding of loop iterations. *Journal of Real-Time Systems*, May 2000.
- [12] C. Healy and D. Whalley. Tighter Timing Predictions by Automatic Detection and Exploitation of Value-Dependent Constraints. In Proc. 5<sup>th</sup> IEEE Real-Time Technology and Applications Symposium (RTAS'99), pages 79–88, June 1999.
- [13] N. Holsti, T. Långbacka, and S. Saarinen. Worst-Case Execution-Time Analysis for Digital Signal Processors. In Proceedings of the EUSIPCO 2000 Conference (X European Signal Processing Conference), September 2000.
- [14] Raimund Kirner and Peter Puschner. Transformation of Path Information for WCET Analysis during Compilation. In Proc. 13<sup>th</sup> Euromicro Conference of Real-Time Systems, (ECRTS'01), June 2001.
- [15] Y-T. S. Li and S. Malik. Performance Analysis of Embedded Software Using Implicit Path Enumeration. In Proc. of the 32:nd Design Automation Conference, pages 456–461, 1995.
- [16] S.-S. Lim, Y. H. Bae, C. T. Jang, B.-D. Rhee, S. L. Min, C. Y. Park, H. Shin, K. Park, and C. S. Ki. An Accurate Worst-Case Timing Analysis for RISC Processors. *IEEE Transactions on Software Engineering*, 21(7):593–604, July 1995.
- [17] S-S. Lim, J. Kim, and S. L. Min. A Worst Case Timing Analysis Technique for Optimized Programs. In Proc. of the fifth International Conference on Real-Time Computing Systems and Applications (RTCSA); Hiroshima, Japan, pages 151–157, Oct 1998.
- [18] G. Ottosson and M. Sjödin. Worst-Case Execution Time Analysis for Modern Hardware Architectures. In Proc. ACM SIGPLAN Workshop on Languages, Compilers and Tools for Real-Time Systems (LCT-RTS'97), June 1997.
- [19] C. Y. Park. Predicting Program Execution Times by Analyzing Static and Dynamic Program Paths. *Real-Time Sys*tems, 5(1):31–62, March 1993.
- [20] P. Puschner and C. Koza. Calculating the Maximum Execution Time of Real-Time Programs. *The Journal of Real-Time Systems*, 1(1):159–176, 1989.
- [21] P. Puschner and A. Schedl. Computing Maximum Task Execution Times with Linear Programming Techniques. Technical report, Technische Universität, Institut für Technische Informatik, Wien, April 1995.
- [22] F. Stappert and P. Altenbernd. Complete Worst-Case Execution Time Analysis of Straight-line Hard Real-Time Programs. *Journal of Systems Architecture*, 46(4):339–355, 2000.
- [23] F. Stappert, A. Ermedahl, and J. Engblom. Efficient Longest Executable Path Search for Programs with Complex Flows and Pipeline Effects. In Proc. 4<sup>th</sup> International Workshop on Compiler and Architecture Support for Embedded Systems, (CASES 2001), November 2001.

# Session II: Tools

Jan Gustafsson

# **Presentations**

## Status of the BOUND-T WCET Tool

Niklas Holsti and Sami Saarinen Space Systems Finland Ltd., Espoo, Finland

The author presented a WCET tool intended for the space industry to begin with, but which now is being marketed in a larger market. The tool analyzes programs on object code level. The author discussed a number of tricky analysis problems, for example due to complex program sequencing logic in certain DSPs.

# WCET Estimation from Object Code Implemented in the PERF Environment

Douglas Renaux, Joa o Goe's and Robson Linhares Laboratory of Embedded Systems Innovation and Technology, Brasil

The authors presented a WCET tool that presents the code in a number of windows. By clicking on or selecting portions of the code, all other windows are updated simultaneously. The tool shows the worst case path. The user had to annotate the code for loop bounds.

The analysis, which is taking place on the object code level, calculates WCET, BCET (Best Case) and TCET (Typical Case).

# Discussion

Q: Has the PERF tool been evaluated outside university?A:No, it hasn't been used outside the development team. It is only a prototype.

**Q**: What about loop annotations ?

A: To work, the tool need loop annotations. These are not checked. We plan to include autmomatic data flow analysis so at least some of the annotations can be calculated automatically.

**Q**: What about infeasible paths?

**A**: If the worst case path is infeasible, the user can (manually) exclude it from the calculation.

**Q**: The results you show are very tight. What type of code do you analyze? **A**: We analyze "real" code which turns out to be simple (functions with 30 - 40 basic blocks, nested loops with number of levels less or equal to 2). Switches are common, and the analysis can be very complicated if they are implemented as tables. **Q**: Are the results safe?

A: We have measured on real hardware, and all our results are safe.

**Q**: What hardware model do you use?

A: It is a computation model, not a mathematical model. There is a problem of correctness.

- Here followed a lengthy discussion on the hardware model used -

**Q**: How to build models that we can rely on? Is there a useful strategy to build correct models?

A: We start with the manual. Then we examine the execution of a number of test functions using an logic analyzer. However, we sometimes need to see what goes on inside the cache (there is such an logic analyzer, too). When processors get more and more complicated, this way of work will not function anymore. On the other hand, the processors used for hard real-time should be simple also in the future.

**Q**: Do you trust the processor manufacturers? **A**: No, we have to verify their numbers.

 $\ensuremath{\mathbf{Q}}\xspace$ : Do you handle data caches?

A: Not yet.

**Q**: How do you handle function calls?

**A**: Each function call is handled independently. Function calls within functions are analyzed first.

**Q**: How are the times included in the Bound-T model?

A: We use the manufacturers numbers. We verify some of them using a simulator, but basically we trust the processor manufacturers.

- The only open processor architecture is the SPARC, since you can get an open desription of what it really does. For other processors, they get more and more complicated, and it gets harder and harder to understand what is really happening inside. We should not accept this state of affairs! We should react!
- We need some type of certification or similar for hardware, with demands for better documentation. Otherwise it shouldn't be allowed to be used in hard real-time systems.

**Q**: How do we get out with our tools into the market?

A: The basic problem seems to be that our results are completely hidden! We need to get out more, for example write articles in magazines outside academia.

- We need an agency of some sort to act as a facilitator. This agency should demand certain standard, certifications or similar. It is necessary to persuade these agencies of the value of WCET tools, not only the users! It is necessary to have some force, since programmers are conservative! The gap is too wide! And there is no-one that is strong enough on his own.
- There are other uses of a WCET tool that can make its way to the market, for example control flow graph analysis. FAA, for example, demands CFG analysis for their software.

- The WCET values does not always have to be very exact. Sometimes less accuracy is enough.
- Maybe it will be most important to be able to prove that the WCET is safe. This is the case for many critical applications. Also it has to be very simple to use, so developers will really use it.
- We need to persuade the managers!
- There is a trend in for example ESA that timing figures are needed. It comes from the top, not from bottom, and comes in the form of standards. Unfortunately, it takes a lot of time.

# WCET Estimation from Object Code implemented in the PERF Environment

Douglas Renaux, João Góes, Robson Linhares Laboratory of Embedded Systems Innovation and Technology – LIT CEFET-PR Brasil

{douglas, goes, robson}@lit.cpdtt.cefetpr.br

#### Abstract

"The estimation of the Worst Case Execution Time of a function produces results that are **safe** and that have a **low error**, even in architectures using pipelines and caches." This is our thesis; in this paper we present results that indicate that this thesis is correct.

The two basic approaches to obtain WCET of a piece of code are estimation and measurement. At LIT, a tool called PERF is under development. This tool uses both approaches to obtain the best of both worlds. Measurement provides precise results, but requires that the target is already built and that the worst possible scenario is being executed, which is often hard to determine. On the other hand, the precision of estimation methods is highly dependent on the complexity of the estimation model. PERF is a design and evaluation environment: a project can be defined, files can be edited, compiled, linked; the resulting code can be analyzed from a timing perspective both via estimations as early as possible in the design cycle. Any tool (commercial or academic) can be inserted in PERF via *plug-ins*. This was the case of the text editor, the compiler and the linker. Hence, PERF is actually a framework to which many tools can be added. PERF works with the object code generated by the integrated tools, in order to obtain execution time limit estimations for functions that compose a real-time systems' software project.

In this paper we present the PERF environment's architecture, with emphasis on the integrated time estimation model and the results obtained using this model.

#### 1 Introduction

Any tool which intends to estimate execution times should take two different domains into consideration: the source code, which the developer usually uses to develop his software project, and the executable code, on which the time estimations are actually performed.

The problem of execution time estimation of a program is usually divided into 3 sub-problems: execution path analysis on source level; source code and machine code correlation; and execution time analysis for each individual machine instruction at each path of the object code.

A strategy based on these 3 sub-problems would search for existing paths in the source code and tries to relate them with the corresponding paths in the object code. Meanwhile, this relationship may not be easily determined, especially in cases where the code optimizations are enabled. The correlation between source-code and object-code may not be trivial, specially for a tool which intends to do this automatically; so, it is possible to concentrate the estimation tool's work on the object-code path's determination and on the estimation of each of its execution times, leaving to the compiler the responsibility of doing the correlations and interpreting the results obtained through the analysis of object code. Another advantage of this strategy comes from the possibility of analyzing object-code present in code libraries, to which the source code is usually not available.

The resolution of the third sub-problem, concerning to the individual execution times of any machine instruction, is affected by the quality of the model that expresses the hardware platform on which the code is to be executed. In order to minimize the estimation error, when compared to the measured time values for a given execution path, this model should consider internal architectural features that have any influence over the execution times of the instructions, such as cache memories and pipelines (if available) [Hennessy 95]. Moreover, it is useful for the model to be reconfigurable, so that the estimation tool is able to address several target architectures of real-time systems. The reconfigurability feature can be available to the tool's user, in a way that allows him to adapt the estimation process to the architectures of his

interest, using configuration parameters; that requires the estimation algorithms to be generic and able to use the configuration parameters of any architecture in a similar way.

### 2 The PERF environment

PERF is a tool, under development at LIT, which intends to be a complete design and evaluation environment. At PERF a software project can be defined, edited, each of its modules can be compiled and linked; moreover, the resulting object code can be analyzed from a timing perspective both via estimation and via measurement, which makes PERF suitable for development of software for real-time systems. In this way, the use of PERF intends to encourage the developer to perform time estimations as early as possible in the design cycle.

PERF's architecture is composed of a central core, controlled through a graphical interface. This core is intended to manage a set of integrated tools, each one performing a determined task in the development process (editors, compilers, linkers, time analyzers, etc.) and configured via a *plug-in*. Any tool (commercial or academic) can be inserted in PERF via *plug-ins*, allowing the developer to use only the functionalities which are strictly needed for the target platform and the type of design /evaluation process of interest.

The PERF environment is a result of a ten year long research in the area of timing analysis; several approaches to determine the execution time of software have been tried and analyzed before the PERF environment was defined. The table below highlights the milestones if this effort.

1992, 1993	The RTX-Parlog language was designed, based on the concurrent
[Renaux 93]	logic programming language Parlog, to provide support for real-time
[Renaux-Dasiewicz 93]	systems. The programming environment of RTX-Parlog included a
	timing tool that was able to determine the execution time of program
	segments (aka basic blocks or program blocks). Here the basic ideas of
	timing analysis were put in practice: defining program segments,
	building control flow graphs where the nodes represent program
	segments, analyzing possible execution paths on a control flow graph
	and combining the execution times of each block.
1994-1995	Use of a logic analyzer to measure the execution time of instrumented
[Renaux 95]	software, i.e. software to which specific instructions were added that
	produced signals on the processor's bus that were detected and stored
	by the logic analyzer
1998	To our built-in-house real-time kernel (PET) we added the capability
[Braga-Renaux 98]	of real-time tracing. The trace is stored in a limited area of the target's
	memory and uses the timer of the kernel as its time base (resolution of
	1 ms). The trace is later downloaded to a PC and analyzed off-line.
1999	PERF – Definition of the architecture
[Renaux-Braga-Kawamura 99]	
1999	TLM – Time-Line Monitor – TLM is a measurement system
[Braga 99]	consisting of a measurement device, physically connected to the
	processor's bus and a visualization software running on a PC. The
	device stores large amounts of trace information and uses a
	nanosecond-resolution clock. After an acquisition the trace is
	downloaded to a PC and analyzed.
1999	FPTE – a tool for static analysis of the execution time of assembly
{Kawamura 99]	code. C and C++ code is analyzed after the compiler generates the
	corresponding assembly code. The major drawback of this approach is
	the analysis of libraries where usually only the object code is
	available.
	PERF analyzes object files
2001 - 2002	PERF's implementation can analyze object code for 80186 CISC
[Gues 01]	incrocontroller and for the powerry 800. The precision of the time
[Linnares 01]	esumation is very good.
[Goes-Linnares-Renaux 01]	
Linnares-Renaux 02	

#### 2.1 Time estimation tool

PERF aggregates an execution time estimation tool, whose estimation process is shown in Figure 1. This figure shows that an important design decision of PERF's estimation tool was to analyze object code, instead of source code or executable code. Object code has all the information that is relevant for timing analysis and it may include debug information that relates the executable code to the source code. Furthermore, by analyzing object code, all the libraries can have their timing information extracted, including commercial libraries.



Figure 1 - Time estimation process used in Perf

Figure 2 shows how PERF presents the different views of a function to the developer. It is worth noticing that the correlation between source and machine code can also be obtained through debug information (as it is shown by the gray shaded areas on the source code, machine code and control flow graph views); nevertheless, the developer can still be asked to provide data flow information necessary for the estimation, such as number of loop iterations, so that the several execution paths found in the control flow graph (CFG) are correctly estimated and the correct values for BCET, TCET and WCET are shown in their respective view.



Figure 2 - The PERF tool presenting several views of a function: source, machine code, control graph and timing information

How does PERF obtain execution time estimations (Figure 1):

#### 1) Analyze the object code and extract the control flow graph.

Each node of the graph represents a program block, i.e., a sequence of machine instructions in which only the last instruction can be a branch (absolute jump, conditional jump, procedure call, software interrupt,...) and only the first instruction is the target of a branch. Control flow graphs are obtained in a function-by-function basis, even for commercial libraries. Each function can have its CFG presented, as well as the possible correlation with the source code (if available).

#### 2) All possible paths in a procedure are analyzed.

Loop structures require information from the user specifying the minimum, typical, and maximum number of iterations. This information can be given in the form of comments in the source code or requested interactively during the analysis (Figure 2).

# 3) For each path, a time estimate (BCET, TCET and WCET) is obtained based on the configured hardware model for the target platform.

The hardware model is the main functionality of the execution time estimation tool. This model addresses the third sub-problem of the execution time estimation problem (the analysis of the individual execution times for each machine instruction) by considering internal architectural features of the target platform's hardware.

The computation model implemented at PERF is divided in two basic parts: the estimation algorithms, which are generic and should work in a similar view, independent of the target platform; and the architecture configuration parameters, obtained from an external *plugin*.

The generic estimation algorithms consider the influence, on machine instruction execution times, of multi-stage pipelines, instruction caches of several sizes and associativity degrees and prefetch queues. It is possible, considering the influence of these features in a correct way, to obtain time estimation values with estimation errors lower than 10%, comparing to measured time values; this is valid not only for RISC architectures but for CISC architectures too. The efficiency of the model's work depends on a correct configuration for the architectural parameters via the configuration plugin; it is conceived as a series of data structures in C++ programming language, which can be configured by a developer using the manufacturer information about the target platform.

For each of the program blocks constituting an execution path, a pre-categorization of the instruction cache behavior for each instruction is performed. This categorization is based on the extended timing schema algorithm, proposed by [Min 94][Lim 94], and adapted in order to address not only the cache but also the pipeline behavior. For this purpose, each instruction's categorization [Mueller 00] is stored in the data structure, not only the total execution time for a program block as it is done in the original algorithm. This approach allows a better precision for the generation of the pipeline information and the prefetch queue behavior, provided that each instruction has its fetch time well determined.

The next step of the estimation process involves the concatenation of the cache categorizations for each of the program blocks of the current path, generating the real cache state for that path. The pipeline information for each path is then generated, using the real cache state and a reservation table for each instruction of the considered architecture. The reservation tables are simplified, as proposed by [Healy 95], in order to describe not only the use of each pipeline stage (for structural conflict determination), but also the register utilization (for data conflict determination).

Each of the instructions of a determined path is concatenated, to generate the pipeline utilization scenario for that path. The concatenation is done according to some rules, which avoid the occurrence of structural and data conflicts, including any bus conflict, which can occur between stages which perform external memory accesses.

The model differentiates the values of BCET and WCET, for each analyzed path, by a pessimistic factor introduced in the algorithms. For this factor to be correctly addressed, the best and worst case execution time values for each machine instruction should be modeled during the configuration of the model; moreover, pessimistic and non-pessimistic considerations are also made, according to the following rules:

- Prefetch queues are optimistic for BCET calculation. Meaning that any instruction that begins its execution (the fetch operation) is immediately removed from the queue; this allows other

instructions to fit into the released space, minimizing the probability of a prefetch contention. On the other hand, the WCET calculation involves the consideration of a pessimistic queue, meaning that the instructions are removed from the queue only when their execution is finished, increasing the possibility of a prefetch contention.

- Bus contentions, due to simultaneous data read/write and prefetch operations, are not considered in the calculation of BCET.

TCET values are obtained through typical-case annotation, provided by the user for the number of loop iterations; in functions where no loop structure is present, the TCET value is the same as the WCET value.

The loop processing, for any analyzed path, is performed using a strategy that eliminates loop redundancies. This decreases the processing times for loops with a great number of iterations and is done according to the following rules:

- the loop is processed and its program blocks are concatenated, generating a macro-block as shown in Figure 3. This macro-block is then concatenated to the timing schema object representing the path's execution until this point; this concatenation represents the first iteration of the loop.



Figure 3 – Loop structure turned into a macro-block

- the second iteration is then processed, through a new concatenation of the macro-block considering the cache and pipeline states after the first iteration.
- The cache and pipeline states after the last and the current iteration are compared. If there is any difference between them, the macro-block is concatenated again and the number of remaining iterations is decreased.
- In case there is no difference, it is presumed that the execution of the further iterations lead to a steady state. It is not needed to concatenate each of the remaining iterations separately, but only add the time of each iteration *n* times to the execution time until this point, where *n* is the number of remaining iterations.

The loop processing of PERF addresses only the cases where any iteration of the loop is executed through the same path, and the existence of multiple paths inside a loop is not completely addressed; for example, in a case where the loop body contains an *if-else* structure and the *then* part is executed in even iterations and the *else* part in odd iterations (Figure 4), the loop must be unrolled and the *n* initial number of iterations should be turned to n/2, each of the iterations considering the test, the execution of the *then* part, the test again and the execution of the *else* part.



Figure 4 – Loop structure with multiple paths

#### 4) A report is generated to inform BCET, TCET and WCET of each function.

The execution time of called functions is included in the execution time of the caller, as long as the execution time of the callee is known. Several iterations may be required to evaluate all these dependencies.

#### 2.2 Hardware model validation

Any processor hardware model, which is configured via a plugin of the time estimation tool, must be validated in order to be used on the time estimation process for real projects.

The configuration process of a model depends on manufacturer's documentation, especially concerning to the instruction set. However, the available documentation is not always precise enough, and that can make necessary the measurement of the execution times of some instructions. Thus, the model must go through a validation process, at which possible inaccuracies can be detected. This step is done by estimating the execution times of test functions, with the typical control structures found on real-time tasks, and comparing those estimated times with times measured on a real platform. If the safety and precision requirements are not satisfied, then the HW model must be revised and tested again.

#### 3 Results

Two processor architectures were initially considered for the time estimation process and modeled, using the proposed hardware model. The first one was a Intel 80C186EC processor, a CISC core with two pseudo-pipelined stages, a 6-byte prefetch queue and no instruction cache memory [Intel 91]; the second architecture was a Motorola PowerPC MPC860 processor, a RISC core with an 8-stage pipeline and a 2-set associative instruction cache memory of 4 Kb [Motorola 98].

The estimation tool was used to estimate the execution time of about a hundred functions, including all the functions of a real-time kernel (PET – [Renaux 96]). For a set of test functions, with simple control structures, the results are shown in Table 1.

The test functions contain loops with a fixed number of iterations, so the measured BCET and WCET have the same value; another set of functions, with variable number of iterations, was also tested and provided similar results. *Test1* and *Test2* contain one and two single loops, respectively, while *Test3* and *Test4* have nested loops with two and three levels, respectively. Although this set of test functions does not have any *if-else* structures, the change of control provided by the execution of the last iteration of each loop and the execution of further instructions is similar to the change of control presented by an *if-else* structure.

80C186EC	Measured time (clock cycles)	Estimated BCET (clock cycles)	Error%	Estimated WCET (clock cycles)	Error%
Test1	421	408	3.09	424	0.71
Test2	972	896	7.82	1005	3.4
Test3	2269	1760	22.43	2286	0.75
Test4	1796	1670	7.02	1812	0.89
MPC860	Measured time (clock cycles)	Estimated BCET (clock cycles)	Error%	Estimated WCET (clock cycles)	Error%
MPC860 Test1	Measured time (clock cycles) 698	Estimated BCET (clock cycles) 674	Error% 3.44	Estimated WCET (clock cycles) 709	Error% 1.58
MPC860 Test1 Test2	Measured time (clock cycles) 698 1208	Estimated BCET (clock cycles) 674 1198	Error% 3.44 0.83	Estimated WCET (clock cycles) 709 1238	Error% 1.58 2.48
MPC860 Test1 Test2 Test3	Measured time (clock cycles) 698 1208 1875	Estimated BCET (clock cycles) 674 1198 1771	Error% 3.44 0.83 5.55	Estimated WCET (clock cycles) 709 1238 1808	Error% 1.58 2.48 -3.57

Table 1- Sample of the results for test functions

The WCET estimation errors for the complete set of functions, compared to measured values, varied from 1% to 7% for the 80C186EC architecture and from 0% to 15% for the PowerPC architecture, depending on whether the instruction cache was enabled or not. Moreover, the results obtained from the estimation tool are all safe (the estimated WCET is never lower than the actual), except for the cases of WCET for *Test3* and BCET for *Test4* estimated for the MPC860 architecture. Safety and precision errors in the PowerPC architecture are due to the inadequate equipment used to validate the model. This is one of the difficulties that can be faced by the developer who intends to configure his own processor model, since the information provided by the manufacturer for each instruction's execution time are not always accurate.

Some weaknesses are still detected on PERF 's execution time estimation tool:

- The inefficiency of solving all kinds of jump tables, used for indirect branch operations (for example, the compilation of most 'switch' structures). PERF is only capable at this moment to find jump tables in architectures for which the base address is encoded in the instruction code (for example, the Intel 80C186EC architecture).
- For optimized code and for code for which the source is not available (for example, libraries), it can be difficult for the developer to determine the exact number of some loops' iterations, when asked by the tool to provide those values.

#### 4 Conclusions and future work

On going work in the PERF environment includes: improving the computational models for the PowerPC architecture to reduce the estimation error, development of computational models of data caches (currently only instruction caches are modeled) and a tool for scheduling analysis is also under development. Also, the problem concerning to jump tables can be solved by allowing the user to interact with the CFG construction; this interaction can include the possibility of creating edges, between program blocks, which could not be automatically determined during the graph analysis.

Some improvements can also be done concerning to the graphical interface and presentation of results to the user. These improvements include the presentation of intermediary steps on the execution time estimation process, allowing the user to visualize the cache and pipeline states at any time.

The results obtained so far with PERF are very encouraging. Our previous work on timing analysis included the analysis of source code, of the assembly listing produced by the compiler and of binary code. Each one has strong aspects, but from our perspective, object code analysis is the best option.

#### 5 References

[Braga-Renaux 98] Braga, André S.; Renaux, D. P. B. "Utilização de Linhas de Tempo para Depuração e Validação Temporal de Sistemas em Tempo Real". In I Workshop de Sistemas em Tempo Real. Rio de Janeiro, RJ, Brazil, May 1998.

- [Braga 99] Braga, André S. "TLM Uma Ferramenta de Apoio ao Teste de Restrições Temporais em Sistemas Dedicados Operando em Tempo Real". Master dissertation. CPGEI, CEFET-PR. Curitiba, PR, Brazil, 1999
- [Braga-Renaux 99] Braga, André S.; Renaux, D. P. B. "Teste de Restrições Temporais Através de Técnicas Funcionais e Estruturais". In II Workshop de Sistemas em Tempo Real p. 25-34. Salvador, BA, Brazil, May 1999.
- [Goes 00] Góes, João Alexandre. Estimação de tempos de Execução de Programas a Partir de Arquivos Objeto". Internal Report T.I. CPGEI, CEFET-PR. Curitiba, PR, Brazil, 2000
- [Goes 01] Góes, João Alexandre. "PERF: Ambiente de Desenvolvimento e Estimação Temporal de Sistemas em Tempo Real". Master dissertation. CPGEI, CEFET-PR. Curitiba, PR, Brazil, July 2001
- [Goes-Linhares-Renaux 01] Góes, João Alexandre; Linhares, R. R.; Renaux, D. P. B. "Estimação de Tempo de Execução de Programas no Ambiente PERF". In Workshop de Tempo Real do XIX Simpósio Brasileiro de Redes de Computadores. Florianópolis, SC, Brazil, May 2001
- [Healy 95] Healy, C. A.; Whalley, D. B.; Harmon, M. "Integrating the Timing Analysis of Pipelining and Instruction Caching". In IEEE Real-Time Systems Symposium, pages 288-297, December 1995.
- [Hennessy 95] Hennessy, J. L.; Patterson, D. A. *Computer Architecture: A Quantitative Approach*. Second Edition. Morgan Kaufmann Publishers, Inc. San Francisco, CA, USA, 1995.
- [Intel 91] 80C186EC / 80C188EC User's Manual. Intel Corporation. Mt. Prospect, IL, USA, 1991.
- [Kawamura 99] Kawamura, Alexandre. "Análise estática de Programas para Estimação de Tempos de Execução". Master dissertation. CPGEI, CEFET-PR. Curitiba, PR, Brazil, 1999.
- [Linhares 01] Linhares, R. R. "Modelamento de Hardware Visando À Estimação do Tempo de Execução de Programas". Master dissertation. CPGEI, CEFET-PR. Curitiba, PR, Brazil, December 2001
- [Linhares-Renaux 02] Linhares, R.; Renaux, D. P. B. "Estimação de Tempo de Execução de Software em Sistemas Em Tempo Real Executando em Plataformas com Pipeline e Cache". In Workshop de Tempo Real do XX Simpósio Brasileiro de Redes de Computadores. Búzios, RJ, Brazil, May 2002.
- [Lim 94] Lim, S-S Et Alli. "An Accurate Worst-Case Timing Analysis for RISC Processors". In IEEE Real-Time Systems Symposium, p. 97-108, December 1994.
- [Min 94] Min, S. L. Et Alli. "An Accurate Instruction Cache Analysis Technique for Real-time Systems". In Proceedings of the Workshop on Architectures for Real-time Applications, April 1994
- [Motorola 98] Power QUICC: MPC860 User's Manual. Motorola, Inc., 1998.
- [Mueller 00] Mueller, F. "Timing Analysis for Instruction Caches". *Real-Time Systems Journal*, vol. 18, number 2/3, p. 209-239, May 2000.
- [Renaux 93] Renaux, D. P. B. "RTX-Parlog: A Concurrent Logic Programming Language for Real-Time Systems". Ph.D. Thesis, Department of Electrical and Computer Engineering, University of Waterloo, Canada. December 1993.
- [Renaux-Dasiewicz 93] Renaux, D. P. B. and Dasiewicz, P. "RTX-Parlog: Real-Time Extended Parlog". In Euromicro '93 Workshop on Real-Time Systems, p. 147-153. IEEE Computer Society Press, June 1993.
- [Renaux 95] Renaux, D. P. B. "Mediç à de tempos de execuç à do software do MMA". Internal Report . LIT, CEFET-PR. Curitiba, PR, Brazil, 1995.
- [Renaux 96] Renaux, D. P. B. "PET A Small Real-Time Support Systems for Microcontrollers without Virtual Memory". Internal Report. LIT, CEFET-PR. Curitiba, PR, Brazil, 1996.
- [Renaux-Braga-Kawamura 99] Renaux, D. P. B.; Braga, André S.; Kawamura, Alexandre. "PERF: Um Ambiente para Avaliaç ão Temporal de Sistemas em Tempo Real". In II Workshop de Sistemas em Tempo Real, p. 76-87. Salvador, BA, Brazil, May 1999.

# Status of the Bound-T WCET Tool

Niklas Holsti and Sami Saarinen Space Systems Finland Ltd Niklas.Holsti@ssf.fi, Sami.Saarinen@ssf.fi

#### Abstract

Bound-T is a tool for static WCET analysis from binary executable code. We describe the general structure of the tool and some specific difficulties met in the analysis of the supported processors, which are the Intel 8051 8-bit microcontrollers, the Analog Devices ADSP-21020 Digital Signal Processor, and the SPARC V7 processor. For the DSP, the problem is the complex program sequencing logic using an instruction pipe-line and nested zero-overhead loops with implicit counters and branches. The solution is to model the full sequencing state in the control-flow graph. For the SPARC, the problems are the register-file overflow and underflow traps, which may occur at calls and returns, and the concurrency of integer and floating-point operations, which may force the Integer Unit to wait when it interacts with the Floating-Point Unit. The traps are modelled with a whole-program analysis. The IU/FPU concurrency is modelled by distributing the potential waiting times onto flow-graph edges in a heuristically optimal way, also using some inter-procedural analysis.

#### **1** Introduction

The Bound-T tool analyses compiled and linked executables to find the WCET, flow graphs, call graphs, and stack usage. Space Systems Finland (SSF) developed the tool with support from the European Space Agency (ESA) for space applications [1][2][3]. SSF is developing the tool further, aiming also at non-space applications.

The target processors currently supported are the Analog Devices ADSP-21020 (a 32-bit floating-point DSP architecture, forerunner of the SHARC), the Intel 8051 (a large family of 8-bit microcontrollers), and the SPARC V7 (a 32-bit RISC general-purpose architecture). All these processors are used in ESA space projects, which is one reason why they were chosen for Bound-T. The other reason was to implement some quite different architectures in order to verify the adaptability of the tool design. Since Bound-T uses only binary code, it is

independent of the source language of the program to be analysed.

The most advanced feature of Bound-T is the automatic analysis of loop-bounds, using a model of the program's arithmetic. The same analysis provides some context-sensitivity by propagating actual parameter values into the analysis of the called subprogram. To supplement the automatic analysis, the user may state assertions on loop bounds, variable ranges, and other useful facts.

The main limitations of Bound-T are currently the lack of analysis of cache memories, a limited analysis of aliasing and dynamic branching, and the difficulty of high-level analysis based on low-level code, especially if the machine word is short, for example a loop with a 16bit counter running on an 8-bit machine. The arithmetic analysis for loop bounds is occasionally very timeconsuming and sensitive to the structure of the program.

The rest of this paper paper is organized as follows. Section 2 discusses the architecture of the tool. Section 3 describes the modelling of the target processor and target program thru abstract data types. Section 4 presents the major analysis phases and methods. Sections 5 and 6 focus on some interesting and difficult problems: the control-flow analysis of the ADSP-21020, and the IU/FPU concurrency in the SPARC. Section 7 reports on the commercialization and section 8 sketches future work.

#### 2 Tool architecture

Bound-T is based on target-specific modules for reading and decoding binary files, generic modules for creating the control-flow graphs and call-graphs, a Presburger Arithmetic package (Omega) [4] for modelling the arithmetic of loop-counters, and an Integer Linear Programming tool ( $lp\_solve$ ) [5] to find the worst-case path.

The architecture of Bound-T was designed to be *adaptable* to different target processors, *extensible* with new kinds and methods of analysis, and *portable* to different host platforms.

The easy part of adaptability is to isolate the targetdependent parts into target-specific modules. The hard part is to make the interface of these modules valid for all
targets. Our approach is to abstract the important aspects of the target processor. This has worked well, as shown by the range of supported targets.

Extensibility is provided in the conventional way by dividing the analysis into phases, with the result of each phase stored in the program-model. This method is limited by the fact that the data structures of the program model are hard-coded (as opposed to a data-base, for example). However, there are hooks to target-specific data and operations which have let us implement the new SPARC analyses in the Bound-T framework.

For portability, we use a portable implementation language (Ada). The current user interface, based on the command-line and text inputs, is trivially portable.

#### 3 Processor and program models

#### 3.1 Processor model

The target processor is modelled by several abstract data types. The most important type is the identifier or *address* of a control-flow *step*. For a simple processor like the Intel 8051, a step-address is just the address of an instruction. For a processor with complex program-sequencing, a step-address can contain much more context (see section 5 for the ADSP-21020 example). The step-address type is the basis for creating the control-flow graphs, where nodes are identified by a step-address.

Another important abstract type models the registers, flags and memories of the processor, or in general any *cell* that can store an integer value. The cell type is the basis for modelling the arithmetic computations and branching conditions. Some cells are just an enumeration, for example all the processor registers that are always statically addressed. Other cells can represent storage addressed in complex or dynamic ways, for example parameters accessed relative to the stack pointer.

These two abstract types divide our model of program state into a *control state* (step-address), modelled by the control-flow graph, and a *data state* (values of cells), modelled by Presburger input-output relations.

Our processor-model is essentially limited to singlethreaded processors, although synchronized internal concurrency is possible. In other words, there may be several functional units running concurrently, as long as they all execute the same instruction stream as in VLIW machines. In the SPARC, the IU and FPU are not that strictly synchronized, and so this processor is in principle out of scope for our model.

#### 3.2 Subprogram model

A subprogram under analysis is represented by a control-flow graph (CFG). A node in the CFG is a basic

block and contains a sequence of steps. A step usually corresponds to a machine instruction, but in special cases instructions may be split into several steps, or consecutive instructions may be bundled into one step. For example, the Intel-8051 may use two consecutive 8-bit immediateload instructions to load an immediate 16-bit value into the 16-bit Data Pointer register, but the arithmetic analysis is easier if the two 8-bit instructions are decoded into one 16-bit load-step in the CFG.

Calls to other subprograms are represented by special CFG nodes (steps) that refer to the callee.

#### 3.3 Timing model

Each step in a CFG has an associated worst-case execution *effort* which depends mainly on the instruction(s) the step represents, but can also depend on the context (via the step address). The effort is a target-specific abstract type, as are the types for the total *work* to execute a sequence of steps and edges, and the processor *power* that determines the number of processor cycles taken by some amount of work. Memory accesses and wait-states are modelled in the effort, work, and power.

Each edge in a CFG has an associated worst-case execution time which typically models two things: (1) the extra time taken to actually branch from a conditional branch instruction, and (2) interference between consecutive instructions. For example, when an ADSP-21020 instruction that uses an Address Generator unit is immediately preceded by a load-register into this address generator, it takes two cycles instead of one cycle. We assign the extra cycle to the edge between the two instructions.

Branch delays due to instruction pipe-lining, where a branch takes effect only after some "delay slot" cycles, are modelled in the step-address (sequencer model), not in the execution time of the branch edge.

Since we use ILP to find the worst-case path in a CFG, we assume that the total execution time of a path is at most the sum of the times for the nodes and edges on the path.

#### 3.4 Arithmetic model

The arithmetic effect of a step is represented by a set of assignments of expressions to cells. The expressions are Presburger formulas, possibly conditional, operating on constants and cell values. Thus, an unconditional formula is a sum of terms where each term is a constant or a cell or the product of a constant and a cell. A conditional formula chooses one of two such sums depending on a Presburger condition, which is a comparison between two Presburger formulas. It is also possible to state that a cell is set to an unknown value. For example, consider a step (an instruction) that increments the register A and sets the Z flag if the result is zero. The effect is represented by the two assignments A' = A + 1 and Z' = if A + 1 = 0 then 1 else 0, where a prime indicates the new value of a cell, so A' = new value of A.

The Presburger formalism could in fact model any Presburger relationship between the "before" and "after" values of the cells, for example A = A' - 1. We use only the functional form (new value = function of old values) to allow other analyses such as *def-use* chaining or constant propagation.

Each CFG edge has a precondition that is a Presburger formula that must be true when this edge is executed. Thus, the precondition is a *necessary* condition for executing the edge, but perhaps not a *sufficient* one. For example, consider a step that decrements register B, jumps to another step if the result is not zero, and otherwise continues to the next step in address order. The edge to the next step has the precondition B = 0 while the edge for the jump has the precondition  $B \neq 0$ . A precondition that is unknown (or too complex to be analysed) is represented by the constant *true*.

Since our aim is only to find loop bounds (not, for example, numerical errors such as division by zero), we only model those instructions and storage-cells that are likely to be used for loop counters. For floating-point instructions we only model the side-effects on the integer computation. On the ADSP-21020 for example, where any 32-bit general register can be used for integer computation or floating-point computation, a floatingpoint computation is considered to yield an unknown register value.

We generally assume that the computation does not overflow or underflow. For the 8-bit Intel-8051 we provide a command-line option to negate this assumption. For example, if register A in the above example is 8 bits wide, unsigned arithmetic is used and overflow is considered, the effect of incrementing A would be encoded as A' = if A = 255 then 0 else A + 1. Unfortunately this creates many conditional assignments which slows down the Presburger solver markedly.

Note that we use a *symbolic* model of the arithmetic; we do not *simulate* the arithmetic by actually computing expressions and assigning their values to simulated cells. The analysis is static and based on solving or simplifying the system of equations and constraints that represents the joint arithmetic effect of all the analysed instructions.

#### 3.5 Program model

The program model consists of all the subprograms under analysis, starting from the "root" subprograms named by the user and including all their callees. We will use the term "call" to mean a specific step, in the CFG of a caller subprogram, that represents a call to a specific callee subprogram.

Each call is associated with a parameter-passing map between the cells in the caller and the cells in the callee. This map is used to propagate bounds on parameter values from the caller to the callee, perhaps for several call levels, in the hope that these parameters define loop bounds. We do not track the other data-flow direction, from callee to caller. The value of any cell that is modified in the callee is considered unknown after the return to the caller.

#### 4 Analysis phases and methods

#### 4.1 Overview

The analysis phases in Bound-T are, in order:

- 1. Reading the target program.
- 2. Instruction decoding and control-flow tracing.
- 3. Arithmetic analysis for dynamic branches.
- 4. Arithmetic analysis for dynamic data accesses.
- 5. Loop bounding analysis.
- 6. ILP analysis to find the worst-case path.

After reading in the binary program and its symbol tables, Bound-T traces the control-flow starting at the root subprogram entry-points. Each instruction is decoded, entered in the CFG as a step (or many steps, or part of a step), and the possible successors (new step addresses) are found and decoded in turn. When a call instruction is found, the callee is added to the set of subprograms to be analysed. This phase terminates when all paths end with a return instruction (a step with no successors) or a dynamic branch (successors so far unknown).

For subprograms with dynamic branches, arithmetic analysis is applied to try to resolve the target addresses of each branch. If this succeeds, the exploration of the control flow is resumed from these addresses. There may be several iterations of flow-analysis and arithmetic analysis.

When the CFGs of all subprograms are complete, arithmetic analysis is applied to try to resolve dynamic data accesses. Unresolved accesses are left as such, and are considered to yield unknown values.

Loop bounding analysis tries to find cells that act as loop counters, to find bounds on the values of these counters, and thus to bound the number of loop iterations.

When loop-bounds in a subprogram are known, we use the Implicit Path Enumeration Technique [6] to find the worst-case path in the subprogram as the solution of an ILP problem where the unknowns are the number of times each CFG node or edge is executed, the constraints are derived from the CFG and the loop-bounds, and the objective is to maximise the total execution time of the subprogram. Subprograms are processed in bottom-up order so that the WCET of each call is known when the caller is processed.

The following sections explain the arithmetic analysis and the loop bounding analysis in more detail. The other phases use conventional methods.

#### 4.2 Arithmetic analysis

In the arithmetic analysis of a subprogram, the arithmetic effects of several steps in the CFG are joined to give the overall effect of some execution path. In mathematical terms, the effect of a step is a relation between the cell values before and after the step, called the *input-output relation*. The input-output relation for a sequence of steps (a path) is computed simply by joining (chaining) the relations of the steps. The preconditions on edges in the path are included as additional constraints in the chain. In an acyclic part of a CFG, a simple one-pass algorithm can compute the input-output relation between the steps, the "incoming" relations to a step where the paths join are combined by set union (disjunction).

Loops are handled by a bottom-up traversal. The body of an innermost loop is acyclic, so we can compute the input-output relation of the body. From this relation, we find the cells that must be loop-invariant (output value = input value). The entire loop is then fused into one step with an effect that approximates the effect of the loop as an input-output relation that keeps the loop-invariant cells constant and does not constrain the other cells, leaving them with unknown values. (To be precise, the relation does include the constraints created by the looptermination paths, from the loop-head to a loop-exit, assuming unknown values at the loop-head.) The next higher (containing) level of loops can then be analysed, fused and approximated in the same way.

In some target processors every instruction sets many flags. This creates many conditional assignments in the effect of the instruction, but most of these are "dead" because the flag is redefined by another instruction before it is used. Before the actual arithmetic analysis as described above, we do a live-variable analysis in the normal way, and include only the live assignments in the input-output relations.

The results of this arithmetic analysis are the inputoutput relations from the subprogram entry-point to each step in the subprogram, or along other interesting paths. From these relations we compute bounds on the values of cells at specific steps, for example bounds on the addresses of dynamic branches and data accesses, or bounds on the actual parameters in calls to lower-level subprograms. Derived or asserted bounds on the parameters of this subprogram, or on local or global variables, can define or sharpen the computed bounds.

#### 4.3 Loop bounding analysis

Loops are bounded with a "syntactic" method, to use Gustafsson's terminology [7]. This method is faster than the abstract interpretation method proposed in [7] but applies only to counted loops.

Each loop is analysed separately as follows, in topdown and flow order. Let the *repeat relation* be the inputoutput relation that represents repetition of the loop, in other words all paths from the loop-head through the loop-body back to the loop-head, with all inner loops approximated as described in the preceding section. All the non-invariant cells in the repeat relation are candidates for loop counters. For each candidate cell, we compute bounds on the candidate's *initial value* before the loop, on the *repeat value* implied by the repeat relation, and on the change in the value implied by the repeat relation. Consider for example this Ada loop:

```
j := 3;
loop
    j := j + 2;
    if ... then j := j + 3; end if;
    exit when j > 9;
end loop;
```

For the cell *j*, the initial value is strongly bounded to j = 3, the repeat value is bounded by  $j \le 9$ , and the change is bounded to  $2 \le \Delta j \le 5$ . Together, these imply that *j* is an up-counter and that the loop-head can be re-entered from the loop-body at most *ceil* (9 - 3 + 1) / 2) - 1 = 3 times, for a total of 4 iterations of the loop.

In this way we can discover up-counters and downcounters. By computing the complement (negation) of the repeat relation we can discover counters that terminate the loop on equality ("exit when j = 9"), but only if  $\Delta j$  is exactly 1.

#### 5 ADSP-21020 program sequencing

Digital Signal Processors often have special support for loops, to speed up vector computation and digital filtering. The ADSP-21020 has an instruction of the form "DO address UNTIL condition" which dynamically creates a loop that starts at the next instruction and ends at the given address. The DO UNTIL instruction sets the processor into a state where fetching the instruction at the given address makes the processor decide whether to repeat or terminate the loop at this address. In other words, from this address control may either continue onwards, or loop back to the instruction after the DO UNTIL, depending on the value of the condition flag two cycles earlier (due to pipe-line lag). Such loops can be nested to six levels and can interact in interesting ways with the delayed branch instructions. For the ADSP-21020, we model in the step-address type the entire three-stage instruction pipe-line (fetch, decode, execute) and the whole six-level loop-nest. Since CFG nodes are labeled by step-address, a CFG edge now represents a transition from one specific pipe-line and looping state to another such state, so the CFG can very precisely model these transitions and their execution time.

An interesting side-effect is that a single instruction can appear as several steps in a CFG. For example, an instruction that follows a conditional delayed branch is decoded twice and represented as two steps, because the step-address for the branch instruction has two successor step-addresses, one that represents the case where the branch will be taken, and the other the case where it will not be taken. These step-addresses differ in the "fetch" stage of the pipe-line model. The delayed branch instruction is thus converted into an immediate branch in the CFG.

For another example, consider a block of instructions that can be entered either through a DO UNTIL instruction, or directly without a DO UNTIL. This whole block is then decoded twice, once in a context with an active loop-level for this DO UNTIL, and once without. However, this case is unlikely to occur in a real program.

#### 6 SPARC IU/FPU concurrency

In the past year and with ESA support, we completed targeting Bound-T to the SPARC V7, including analysis of the register-file overflow and underflow traps and of the concurrency between the Integer Unit (IU) and the Floating Point Unit (FPU). In this paper, we focus on the IU/FPU concurrency.

#### 6.1 The problem

In the SPARC V7, the IU is responsible for fetching all instructions and for executing integer instructions. Each floating-point (FP) instruction is forwarded to the FPU which executes the instruction while the IU fetches and executes new integer instructions. When a new FP instruction is found, but the FPU is still busy, the IU must wait until the FPU has finished the first FP instruction. The delay depends on the execution time of the first FP instruction, which increases the delay, and on the amount of IU computation between the two FP instructions, which decreases the delay. In the worst case, the delay can be nearly 80 cycles. An integer instruction is typically executed in one or two cycles.

One difficult aspect of this problem is that the delay depends on the path taken by the IU between the two FP instructions; this can be a large number (up to 80) of integer instructions. The delay is not an interaction between two *consecutive* instructions. Another difficult aspect is that the integer-execution time appears with a negative sign in the expression for the delay, which means that we would need the *best*-case IU time for computing the *worst*-case delay.

#### 6.2 Finding delayed paths

The first phase in the IU/FPU analysis is thus to find the (potentially) *delayed paths* between two FP instructions (or from an FP instruction to itself), where the intervening integer computation is too brief to ensure that the first FP instruction has finished before the second one should start. Delayed paths are found simply by depthfirst CFG traversals starting at each FP instruction and propagating the maximum remaining FPU-busy time along all paths, until another FP instruction is found or the FPU is sure to have finished the first FP instruction.

We avoid the need for best-case IU times by a principle we call "hurry up and wait". If the IU executes the path between the two FP instructions *faster* than the worst-case bound, it will just have to wait *longer* for the FPU to finish the first FP instruction. In fact, the WCET of the first FP instruction is a worst-case estimate for the whole path from the start of the first FP instruction to the start of the second FP instruction, including the FPU-busy delay before starting the second FP instruction. Thus, we phrase the analysis in terms of the total time for each delayed path, not in terms of the actual delays.

#### 6.3 Assigning delays to edges

Our objective is to minimize the pessimism, so that the FPU-busy delay is associated only or mainly with the paths that really incur delay in execution. The simplest approach would be to assign the worst-case delay time to each edge that enters the second FP instruction and that is part of a delayed path. This is pessimistic if this edge is part of the overall worst-case path and the worst-case path is not in fact delayed at this FP instruction, for example because it does not execute the first FP instruction. It could be less pessimistic to assign the delay to an edge that leaves the first FP instruction, if this edge is used only by the delayed path.

To avoid such pessimism, we use an ILP approach to distribute the worst-case blocking delay from *all* delayed paths onto *all* the edges in delayed paths so that the added pessimism is minimized (using a heuristic goal function, however). In this ILP problem, the unknowns are the additional delays to be assigned to the edges on delayed paths; the constraints are that the total execution time (including these additional delays) of every delayed path must be at least the WCET of the FP instruction at the start of the path; and the objective is to minimize an expression that estimates the pessimism.

#### 6.4 Is it cheating?

When the FPU-busy delays have been distributed to the edges of delayed paths, the subprogram timing model — the CFG with times assigned to nodes and edges — is no longer a worst-case model, since the time assigned to the last edge and node of a delayed path is not necessarily an upper bound on the actual execution time of this edge and node. However, we can prove that the normal Bound-T WCET analysis (which assumes a worst-case model) still gives an upper bound on the execution of the whole subprogram, thanks to the constraints imposed on the distribution of the delays.

#### 6.5 Inter-procedural aspects

To handle FPU-busy delays between FP operation pairs that cross subprograms, a bottom-up interprocedural analysis assigns each subprogram a minimum *margin* of pure IU execution time on entry (before the first FP instruction is reached), and a maximum *legacy* of remaining FPU execution time on return (due to the last started FP instruction). These summary values are associated with the call steps in the analysis of the calling subprograms.

#### 7 Marketing and commercialization

To commercialize Bound-T we have contacted a number of potential development partners, tool distributors and tool users, but the progress is very slow. WCET analysis is still poorly known, and it is often hard to make people understand what it does and how it differs from debugging, simulation and testing.

Partly because of the marketing delays, and partly because of staff shortages, no technical development is currently under way. SSF will itself use the SPARC V7 version in a major project that develops the on-board platform software for the ESA GOCE satellite (Gravity and Ocean Circulation Explorer). This will test the specific SPARC analysis methods described above, as well as the generic abilities of the tool.

We still hope to make Bound-T a commercially available tool, but the near-term plans are vague and depend on finding partners or users. We will gladly supply evaluation copies of Bound-T. The host platforms are Sun Solaris, Intel Linux, and Intel Windows (using CygWin and a command-line interface).

#### 8 Future work

It is evident that the range of target processors should be increased and updated. Candidates for new targets include ARM, MIPS, AVR, and other microcontrollers. We also have plans for several generic technical improvements. The derived loop-bounds could be used to sharpen the resolution of dynamic data accesses and dynamic branching, and the latter two should be iterated when needed, since some dynamic branching depends on dynamic data addressing (switch tables).

In the loop-bound analysis, it would be better to compute bounds on the difference between the initial value and the repeat value of a counter cell, instead of separate bounds on the two values. This would let us bound loops of the form "for j in  $n \dots n + 10$ " even when the value of n cannot be bounded statically.

Nested loops where the bounds of the inner loop depend on the counter of the outer loop will currently yield pessimistic WCETs, because the worst-case bounds on the inner loop are assumed to hold for all executions of the outer loop. We don't have a clear idea how this analysis could be improved in the Presburger method.

Better aliasing analysis and optional levels of aliasing analysis will probably be necessary for large target programs. Finally, while the current command-line interface is quite workable, a GUI would be convenient for browsing the analysis results of large programs.

#### **9** References

[1] ESTEC/Contract No. 13362/77/NL/FM, "DSP Execution Time Estimation".

[2] N. Holsti, T. Långbacka and S. Saarinen, "Worst-Case Execution Time Analysis for Digital Signal Processors", *X European Signal Processing Conference*, EUSIPCO 2000.

[3] N. Holsti, T. Långbacka and S. Saarinen, "Using a Worst-Case Execution Time Tool for Real-Time Verification of the DEBIE Software", *Proceedings of the DASIA 2000 (Data Systems in Aerospace) Conference* (ESA SP-457, ISBN 92-9092-669-4, September 2000), pp. 307-312.

[4] W. Pugh et. al., The Omega Project: Frameworks and Algorithms for the Analysis and Transformation of Scientific Programs, University of Maryland, http://www.cs.umd.edu/projects/omega.

[5] M. Berkelaar, ftp://ftp.ics.ele.tue.nl/pub/lp\_solve.

[6] Y-T. S. Li and S. Malik, "Performance Analysis of Embedded Software Using Implicit Path Enumeration". In *Proc. of the 32:nd ACM IEEE Design Automation Conference (DAC95)*, 1995, pp. 456-461.

[7] J. Gustafsson, Analyzing Execution-Time of Object-Oriented Programs Using Abstract Interpretation, Uppsala University (Ph.D. Thesis, DoCS 00/115, ISSN 0283-0574) and Mälardalen University (MRTC 00/10, ISSN 1404-3041), May 200.

## Session III: Industrial views

Peter Puschner

### Presentations

This session started with three talks on very complementary views of WCET analysis.

In the first presentation "You Can't Control what you Can't Measure, OR Why it is Close to Impossible to Guarantee Real-Time Software Performance on a CPU with Onchip Cache" Nat Hillary discussed different strategies to measure the performance of software. He stated that the measurement strategy to be used for assessing the performance of a piece of code depends on the application and its timing requirements, in particular the temporal accuracy demanded by the application. Depending on the application needs the choice of analysis tools mentioned ranges from logic analyzers to in-circuit emulators, hardware-assisted software performance monitors, and software-assisted software performance profilers. He argued that cycle-accurate measurements are really not necessary for every application.

Second was a talk by Christian Ferdinand about "Validation by Static Analysis and Abstract Testing" for software to be executed on high-performance processors like the Coldfire 5307 or the PowerPC 755. He mentioned that it is mainly the cache that makes WCET analysis difficult. Switching off the cache in order to make the prediction easier is not an option as this would slow down computations by a factor of 30. The analyses get even more complicated due to timing anomalies, e.g., due to prefetching a cache miss on a memory access does not necessarily cause a worst-case scenario wrt. execution time. The second part of the talk comprised a WCET tool description and an illustrative demo of the tool.

In the third presentation Tullio Vardanega summarized the involvement and interest of ESA in WCET analysis. ESA considers WCET analysis as part of their efforts in building reliable real-time systems for space applications. WCET analysis is considered of one of the important issues besides the specification and design, scheduling analysis and testing of real-time software. To achieve their aim, ESA funds the development of practices and prototype tools. The results from these project are evaluated and find their way into recommendations issued by ESA.

#### Discussion

The following discussion centered around the static analysis vs. measurements issue. It was repeatedly stated that the two techniques complement each other and that both must be used together. Neither testing (measurements) nor static analysis alone are sufficient to get certainty that a piece of code meets its deadlines under all circumstances.

In the early phases of software development measurements are considered as a simple way to get a coarse idea about the performance of a section of code. In the later phases of the software development cycles, measurements are important to validate the results of the static analysis on the real target. Static analysis is seen to be a means to obtain WCET estimates early in the software development cycles, when the target system is possibly not yet available or the application software is not yet ready to run in its entirety. Later on, the detailed information about the flow facts of the worst-case paths give important hints for the generation of input data for the measurements. These measurements are then, in turn, used to validate the static analysis.

In the future, on-chip trace facilities of modern processors will allow programmers and software testers to get more information out of measurement runs than was possible in the past. On-chip real-time trace units with a small memory (e.g., 8 to 16Kbyte on the ARM) will log important event and trace information about execution paths. This information can be downloaded for further evaluation once a measurement series has been completed, thus reveiling the execution details of the observed executions.

### You Can't Control what you Can't Measure, OR Why it's Close to Impossible to Guarantee Real-time Software Performance on a CPU with on-chip cache

Nat Hillary Manager of Technical Marketing Applied Microsystems Corp. <u>nath@amc.com</u> Ken Madsen Manager, Product Marketing Wind River Systems, Inc. ken.madsen@windriver.com

June 3, 2002.

#### 1. Abstract

Steady increases in CPU core speeds continue to extend the range of applications for computer-based solutions, resulting in the creation of ever more responsive systems. At these higher core speeds, on-chip cache architectures are used to prevent the CPU from stalling when accessing relatively slow off-chip memory. In normal operation, most fetch-execute cycles occur internally, guaranteeing the execution of the maximum instructions per second. However, this also serves to hide the state of executing code from the user. Given the fact that it is not possible to directly monitor the execution of code within such a CPU when running at full speed, is it possible to guarantee and control the performance of Real-Time software on cache-based CPU architectures?

This paper investigates this issue by first offering a definition of Real-Time software, together with a discussion on what must be measured to prove that the system will meet its performance objectives in all circumstances. The range of currently available software performance monitoring technologies and techniques currently available will be discussed, together with a summary of the pros and cons of each measurement technique.

As with all measurements in science, it is impossible to measure the execution time of Real-Time software without affecting the system. Nevertheless, a range of technologies and techniques are available for monitoring the execution speed of Real-Time software, ensuring that software performance deadlines when executing on a CPU utilizing on-chip cache can be achieved and controlled.

#### 2. Real-Time Software

Real-Time software is simply code for which the time at which the output is produced is significant. This is usually because the input corresponds to some movement in the physical world, and the output has to relate to that same movement. The lag (delay) from the input time to output time must be sufficiently small for acceptable timeliness.

Because Real-Time software has performance criteria included in its specifications, it is essential that software execution performance be monitored at every step during development, from the writing of Interrupt Service Routines (interrupt service routines) to time-critical sections of application code. So what technologies and techniques may be used to measure software execution performance, and what are the implications of using them with a cache-based CPU?

Starting from board bring-up, the technologies most commonly employed for measuring software execution speeds are:

- □ Logic Analyzers
- □ In Circuit Emulators
- □ Hardware-assisted software performance monitors
- □ Software-assisted software performance profilers In general, these technologies are applied to Real-

Time software performance monitoring using one or more of the following fundamental measurement methods:

- Determining where the system is spending its time (e.g. profiling)
- Monitoring the ability of critical sections of code to meet their deadlines
- □ Measuring a systems response to external events.

#### 3. Logic analyzers

Typically used to monitor multiple digital hardware signals simultaneously, logic analyzers may also be used to make high-resolution software performance measurements, normally for measuring the systems response to external events, or for monitoring the execution speed of critical sections of code. Although it is not their forte, they may also be used for performance profiling.

For CPUs utilizing on-chip cache, these types of software performance measurements, when made with logic analyzers, require external CPU signal lines to be asserted when particular lines of code are reached. This results in very high-resolution timing measurements.

An example of measuring the systems response to external events is monitoring interrupt latency times. The technique for this is fairly straightforward; a single command is placed at the entry to the ISR that asserts a signal on an external CPU pin (e.g. a spare chip select or programmable I/O pin, which will not stall the CPU). The logic analyzer is then used to measure the interval between an interrupt occurring and the CPU signal marking entry to the software routine being asserted.

A typical technique for making high-resolution timing measurements of critical sections of code requires that a signal assertion command be inserted at the entry and exit points of the critical section of code. The Logic Analyzer is then used to measure the interval between the two signal assertions. Modern Logic Analyzers (such as the TLA series from Tektronix) extend this technique, allowing specific networking signals (such as Ethernet packets or ATM cell contents) to be used as hardware trigger events.

By using more extensive instrumentation (e.g. adding assertion statements to salient points in code such as function entries, exits, branch points, etc.), logic analyzers may also be used for performance profiling.

For most applications, a prohibitive number of unique off-chip signals are required in order to correctly identify each unique point in code. In this case, external memory writes may be used to ensure that enough unique instrumentation points are used for the measurements to be meaningful. However, as this profiling technique requires external memory writes, it is more intrusive than the techniques described above. It is therefore not recommended that this profiling method be used for making the type of deadline measurements described earlier for applications with extremely tight deadlines.

Logic Analyzers typically do not gather performance data over a statistically long period. It is therefore

necessary to use analytical techniques to ensure that the correct conditions are created so that particular measurements accurately reflect the worst-case execution time of a particular section of code.

In some rare cases, inserting additional code into an application degrades the performance to a point where the system's Real-Time characteristics are not being met. In this case, 'black box' performance testing techniques are required, where measurements are made at points external to the CPU. E.g. the response time between a particular Ethernet packet arriving and the system responding might be measured using a Tektronix TLA Logic Analyzer.

#### 4. In circuit emulators

Typically used in the early debug stages of target board bring-up, In Circuit Emulators (ICEs) may also be used for software performance measurements. Traditionally, the Real-Time bus trace capability was the most significant feature of an ICE for non cache-based CPUs. Real-Time bus trace may be used for measuring the systems response to external events, or for monitoring the execution speed of critical sections of code.

Asides from the lack of profiling data, this Real-Time bus trace is the ideal solution for performance monitoring of true Real-Time software. However, it requires an offchip fetch-execute cycle to occur in order to monitor what's going on.

Modern cache-based CPUs tend not to have full ICE solutions available. Instead, CPU serial Test Access Points (TAPs) are used for processor emulation control.

Most TAP emulation solutions do not have Trace measurements. Triggering timing measurements with a TAP emulator requires the use of hardware or software breakpoints, which are intrusive. In addition, serial TAP buses are slow (typically 33 MHz) by comparison to processor speed and events are detected asynchronously to their occurrence. Any timing measurements made via this bus are going to be subject to inaccuracies; monitoring the execution speed of a 400 MHz CPU core by sending information through a significantly slower serialized communications bus is not an ideal solution.

Traditionally the ideal solution for making software performance measurements on the fly, contemporary ICE solutions rarely support the features required to make deterministic timing measurements of code.

#### 5. Hardware-Assisted Software Performance Monitors

An extension of in circuit emulation technology, hardware assisted software performance monitors, such as the CodeTEST product from Applied Microsystems, are designed specifically to measure software performance.

This technology requires the combination of software instrumentation and hardware data collection. It may be used to monitor low-level code (such as interrupt service routines), application level code and also RTOS activity. In addition, time stamping may be triggered by external events, making the timing of hardware/software interactions (such as interrupt latencies) possible.

The source code instrumentation technology is used during compilation to add tags to salient points in code. Each instrumentation point equates to a single 32 bit offchip write to memory. This introduces the same profiling inaccuracies as with the Logic Analyzers above, so it is not recommended that this level of instrumentation be used for measuring the performance of critical sections of code.

As with the Logic Analyzer solution, monitoring the performance of critical sections of code requires that an instrumentation tag be placed only at the entry and exit points of the critical section of code (in some instances, such as an infinite loop forming the basis of a task, a single instrumentation point may be used). During code execution, the hardware data collection agent makes highresolution measurements of the time spent between the two instrumentation points.

Any measurements made using this technology may be gathered over a significant period of time, with the automatic collation of minimum, maximum and average execution times. When used to measure the performance of critical sections of code, the overhead of each off-chip write is minimal and easy to calculate, making the measurements that this technique provides highly accurate and deterministic.

This technology also provides the best method for general code optimization, by providing application level profiling data that identifies where the system is spending its time, ensuring that optimization efforts are focused on the right areas.

The 'call-pair' data provided by this technology may also be used to improve software performance. 'Callpairs' measurements identify highly inter-dependent functions that make good candidates for either inlining, fixing in cache, or being located close to one another in the link map of the application (increasing the probability

of highly interdependent functions being simultaneously co-located in cache).

#### 6. Software-Assisted Software Performance Profilers

Worthy of mention because of their dominance in the desktop marketplace, software-assisted performance profilers use a variety of techniques for monitoring where an application is spending its time. If this technology is ever used during the development of a Real-Time system, it is used to aide optimization efforts, and not to measure any of the Real-Time characteristics of the code.

Typically consisting of an in-target data collection agent and either code instrumentation or stack/IP sampling, the potential of this technology is intriguing for two reasons. First, these techniques do not require any off-chip accesses in order to make their measurements. Secondly, solutions based on these techniques tend to be extremely easy to use.

On the other hand, these techniques rely on a target based data collection agent, which is intrusive. Anv techniques based on stack/IP sampling are also prone to aliasing, and in require higher levels of intrusion to improve their accuracy.

#### 7. What Level of measurement accuracy is required?

For Real-Time systems, 'Real-Time' does not necessarily equate to 'real-fast'. The environment in which a system must operate dictates the performance criteria of Real-Time software. A pacemaker, for instance, must respond to specific physiological events within a specific time period before permanent damage to the heart ensues (response times in the 100's of mS). Meanwhile, a commercial flight control system must process and respond to thousands of inputs a second, from pilot commands to air data (response times in the mS).

With modern CPUs capable of processing in excess of 2 billion instructions per second, is it really necessary to measure software performance on a per instruction basis? The simple answer is no, provided that:

- Worst-case response/execution times of a system are monitored, verified and managed
- Enough information is to hand during software creation to ensure that the system performance objectives can be met.

From this, then, the question then arises whether this is achievable with CPUs utilizing on-chip cache.

For extremely high accuracy software performance measurements of worst-case execution time (e.g. nS accuracy), Logic Analyzers must be used. Alternatively, if uS accuracy of software performance is required, then hardware assisted software performance monitoring technologies show the most promise. The only question is whether the performance impact of the off-chip writes that these technologies require is prohibitive, or not. This is worth a more detailed consideration.

When measuring the worst-case execution time of a critical section of code (e.g. the main loop in a control function) using this technology, a single write statement is required. Timing is started when the write occurs the first time, and then the interval to the next occurrence is timed. But what overhead does this introduce?

Consider a typical environment where a target system is using a 100 MHz external CPU bus that requires 3 clock cycles to complete a write operation. In this instance, the delay imposed by each write operation would be a deterministic 30 nS.

The impact of a 30nS delay per cycle in the time critical code of a Real-Time system is negligible. The impact of being able to deterministically measure the worst-case execution time of the software under development with uS accuracy, however, is not. This lends great credence to the power of hardware assisted software performance monitoring technologies, especially when these technologies may be used to gather timing information on the critical sections of code over a significant period of time, ensuring the true worst-case execution time is understood.

#### 8. Conclusion

It is an age-old dilemma in science; how can you measure something without affecting it? When it comes to measuring the performance of Real-Time software, the simple answer to this is - you can't. Add a CPU that utilizes on-chip cache, and the situation only gets worse. It is imperative, therefore, that the right performance measurement technique be used for the software being created. If the Real-Time nature of the software under development requires a timing accuracy in the nS range, then a Logic Analyzer must be used for software performance measurements. It must be understood. however, that data can only be gathered over a limited measurement period. Therefore, careful consideration must be made in the creation of the stimulus or circumstances to make sure that the worst case scenarios are represented for measurement and analysis.

Traditionally, Logic Analyzers required intimate knowledge of memory implementations on the target hardware, thus they provided very little functionality for software engineers. However, new products such as LA Trace from Wind River Systems abstracts the bus implementation from the user making it easy for software engineers to configure the circuitry of a Logic Analyzer to make complex timing measurements. Furthermore, Wind River's LA Trace is able to leverage RTOS knowledge to present acquired information relative to RTOS threads and events.

On the other hand, if you want information in the uS range, use the type of hardware assisted software performance monitoring technology available with the CodeTEST product from Applied Microsystems. This not only provides accurate one-shot timing information, but it also gathers performance information over an indefinite period of time, ensuring that the worst-case execution time of the software being measured is encountered. In addition, the same technology provides function profiling data that greatly enhances optimization efforts, and call-pair information that enables immediate performance improvements through in-lining or prudent link-map ordering.

Real-Time bus trace data from in circuit emulators have traditionally the fall back solution for Real-Time software performance measurements. Most modern CPUs utilizing on-chip cache, however, only have serial Test Access Point emulation solutions without Real-Time bus trace capabilities. Emulators do not, therefore, provide the performance information that they once did.

Software only profiling solutions, popular in the desktop market, are too intrusive and/or inaccurate to make accurate worst-case execution time measurements for Real-Time systems. However, they do provide the profiling information that may be used to yield significant performance improvements during code optimization.

As with all measurements in science, it is impossible to measure the worst-case execution time of Real-Time software without affecting the system. Nevertheless, technologies are available that are appropriate for the required level of accuracy, ensuring that the Real-Time nature of software executing on a CPU utilizing on-chip cache can be controlled.

#### 9. References

[Clement02] Marc Clement, Integrating & Expanding Embedded Design Toolset, Tektronix White Paper

[DeMarco98] Tom Demarco, Controlling Software Projects: Management, Measurement, and Estimates, Prentice Hall PTR/Sun Microsystems Press, 1998

[Grehan98] Rick Grehan, Robert Moote, Ingo Cyliax, Real-time Programming: a guide to 32-bit embedded development, Addison Wesley, 1998

[Heath98] Steve Heath, Embedded Systems Design, Newnes, 1998

[Hillary01] Nat Hillary, Guaranteeing the Performance of Real-Time Systems, Applied Microsystems White Paper, 2001

[Hillary02] Nat Hillary, Gaining Control of Software Performance, Applied Microsystems White Paper, 2001

[Smith98] Connie U. Smith, Lloyd G. Williams, Software Performance Engineering for Object Oriented Systems: A Use Case Approach, Performance Engineering Services, 1998

[Wettersten96] Erik Wettersen, Implementing Performance Engineering, Presentation from InterWorks Technical Users Forum of Interex Conference April 21-24, 1996 Marc Langenbach \* † mlangen@cs.uni-sb.de Christian Ferdinand <sup>†</sup> ferdinand@absint.com Reinhard Wilhelm \*

wilhelm@cs.uni-sb.de

#### Abstract

Computers controlling potentially hazardous machinery and systems are expected to always execute in time. Consequently, for the modeling and planning of embedded systems it is essential that the worst case execution time (WCET) of all program tasks is known.

Modern processor components like caches and pipelines complicate the task of determining the WCET considerably, since the execution time of a single instruction may depend on the execution history. The safe yet almost never valid assumption of a cache hit never occurring results in a serious overestimate of the WCET. Overestimates of the WCET in turn result in an overscaled hardware design.

Our approach to WCET prediction includes a static prediction of cache and pipeline behavior, enabling a much better upper limit to be computed for the WCET. The result is that safety-critical systems can be designed smaller and more cost-effectively.

#### 1 Introduction

Hard real-time systems have specified deadlines for their tasks. It is the duty of the developer to guarantee that the tasks making up the system will always meet these specified deadlines. When it comes to processors with caches and complex pipelines, computing sharp upper bound on the worst-case execution time (WCET) is of critical importance.

There is a tremendous gap between the cycle times of modern microprocessors and the access times of main memory. Caches are used to overcome this gap in virtually all performance-oriented processors (including highperformance microcontrollers and DSPs). Pipelines enable acceleration by overlapping the executions of different instructions. The consequence is that the execution behavior of the instructions cannot be analyzed separately since it depends on the execution history.

Cache memories usually work very well, but under some circumstances minimal changes in the program code or program input may lead to dramatic changes in cache behavior. For (hard) real-time systems, this is undesirable and possibly even hazardous. The widely used classical methods of predicting execution times are not generally applicable. Software monitoring or the dual loop benchmark change the code, what in turn has impact on the cache behavior. Hardware simulation, emulation, or direct measurement with logic analyzers can only determine the execution time for one input. This can not be used to infer the cache behavior for all possible inputs in general. Making the safe—yet for the most part—unrealistic assumption that all memory references result in cache misses results in the execution time being overestimated by several hundred percent.

#### 2 Timing Validation

Real-time systems are typically composed of a set of tasks with specified deadlines (mostly dictated by the surrounding physical environment). A schedulability analysis has to be performed in order to guarantee that all timing constraints will be met (*timing validation*). All existing techniques for schedulability analysis require the worstcase execution time of each task in the system to be known prior to its execution. Since this is not computable in general, estimates of the WCET have to be calculated. These estimates have to be safe, i. e., they must never underestimate the real execution time. Furthermore, they should be tight, i. e., the overestimate should be as small as possible.

#### **3** WCET Computation

The determination of the WCET of a program task is composed of several different tasks:

- Value Analysis: computation of address ranges for instructions accessing memory
- Cache Analysis: classification of memory references as cache misses or hits

<sup>\*</sup>FR Informatik, Universitaet des Saarlandes, Postfach 15 11 50, 66041 Saarbruecken

<sup>&</sup>lt;sup>†</sup>AbsInt Angewandte Informatik GmbH, Stuhlsatzenhausweg 69, 66123 Saarbruecken



Figure 1. The structure of the analysis. The AIP file contains parameters for the analyses, the PER file the results obtained from the combined cache/pipeline analysis.

- **Pipeline Analysis:** prediction of the behavior of the program on the processor pipeline
- **Path Analysis:** the determination of a worst-case execution path of the program.

Many of these tasks are quite complex for modern microprocessors and DSPs.

The arrangement of the tasks is described in Figure 1. The results of the value analysis are used by the cache analysis to predict the behavior of the (data) cache. The results of the cache analysis are fed into the pipeline analysis allowing the prediction of pipeline stalls due to cache misses. The combined results of the cache and pipeline analyses are used to compute the execution times of program paths. The separation of WCET determination into several phases has the additional effect that different methods tailored to the subtasks can be used. In our case, the value analysis, the cache analysis, and the pipeline analysis are done by *abstract interpretation* [3], a semantics-based method for static program analysis. Path analysis is done by integer linear programming.

#### 3.1 Reconstruction of The Control Flow from Binary Programs

The starting point of our analysis framework (see Figure 1) is a binary program and additional user-provided information about numbers of loop iterations, upper bounds for recursion, etc.

In the first step a parser reads the compiler output and reconstructs the control flow [18, 19]. This requires some

knowledge about the underlying hardware, e. g., which instructions represent branches or calls. The reconstructed control flow is annotated with the information needed by subsequent analyses and then translated into CRL (Control Flow Representation Language)<sup>1</sup>. This annotated control flow graph serves as the input for microarchitecture analyses.

#### 3.2 Value Analysis

The value analysis determines ranges for values in registers and by this it can resolve indirect accesses to memory.

#### 3.3 Cache Analysis

The cache analysis classifies the accesses to main memory, i. e., whether or not the needed data resides in the cache. The following categories are used:

- **always hit (ah)** The memory reference will always result in a cache hit.
- always miss (am) The memory reference will always result in a cache miss.
- **persistent (ps)** The referenced memory block will be loaded at most once
- **not classified (nc)** The memory reference could neither be classified as **ah** nor **am**.

The cache analysis used here is based upon [4].

 $<sup>^1{\</sup>rm CRL}$  is a human-readable intermediate format designed to simplify analyses and optimizations at the executable/assembly level.



Figure 2. The MCF 5307 Pipeline

#### 3.4 Pipeline Analysis

The pipeline analysis models the pipeline behavior to determine execution times for a sequential flow (basic block) of instructions, as done in [16, 17]. It takes into account the current pipeline state(s), in particular the resource occupancies, the contents of prefetch queues, the grouping of instructions, and the classification of memory references as cache hits or misses. The result is an execution time for each instruction in each distinguished execution context.

#### 3.5 Path Analysis

Following from the results of the microarchitecture analyses, the path analysis determines a safe estimate of the WCET. The program's control-flow is modeled by an integer linear program [21, 20], so that the solution to the objective function is the predicted worst-case execution time for the input program. A special mapping of variable names to basic blocks in the integer linear program enables execution and traversal counts for every basic block and edge to be computed.

#### 3.6 Analysis of Loops and Recursive Procedures

Loops and recursive procedures are of special interest, since programs spend most of their runtime there. Treating



Figure 3. Map of formal pipeline model

them naively when analyzing programs for their cache and pipeline behavior will result in a high loss of precision.

The following observation can be made frequently: the first execution of the loop body usually loads the cache and subsequent executions find most of their referenced memory blocks in the cache. Hence, the first iteration of the loop often encounters cache contents quite different from those of later iterations. This has to be taken into account when analyzing the behavior of a loop on the cache. A naive analysis would combine the abstract cache states from the entry to the loop and from the return from the loop body, thereby losing most of the contents. Therefore, it is useful to distinguish the first iteration of loops from the others.

A method has been designed and implemented in the program analyzer generator PAG [1], which virtually unrolls loops, the so-called VIVU<sup>2</sup> approach. Memory references are now considered in different execution contexts, essentially nestings of first and non-first iterations of loops.

#### 4 WCET Analyzer for the MCF5307

The ColdFire family of microcontrollers is the successor of the well known M68k architecture of Motorola. The ColdFire 5307 [11] is an implementation of the Version 3 ColdFire architecture. It contains an on-chip 4K SRAM and a unified 8K data/instruction cache.

<sup>&</sup>lt;sup>2</sup>Virtual inlining, virtual unrolling

It implements a subset of the M68K opcodes, restricting opcodes to two, four, or six bytes, thereby simplifying the decoding hardware. The CPU core and the external memory bus can be clocked with different speeds (e. g. 20MHz bus clock and 60MHz internal core clock).

The MCF5307 has two pipelines decoupled by a instruction buffer (see Figure 2): a fetch pipeline fetches instructions from memory, partially decodes them, performs branch prediction, and places the instructions into an instruction buffer, consisting of a FIFO with eight entries (complete instructions). The execution pipeline consists of two stages that obtains complete instructions from the instruction buffer, decodes and executes them.

The WCET analyzer [5] consists of a value analysis, an integrated cache and pipeline analysis and a path analysis.

The implementation of the pipeline analysis is based on a formal model of the pipeline (see Figure 3). In this formal model, a concrete pipeline state consists of several *units* with inner *states* that communicate with one another and the memory via *signals*, and evolve cycle-wise according to their inner state and the signals received.

The decomposition into units accounts for reduced complexity and easier validation of the model. Units often map directly to pipeline stages, but also may represent more than one stage or introduce virtual pipeline stages that are not present in hardware but facilitate the design of the pipeline model (cf. the store stall timer).

Signals may be *instantaneous*, meaning that they are received in the same cycle as they are sent, or *delayed*, meaning that they are received one cycle after they have been sent. Signals may carry data with them, e.g. a fetch address. Note that these signals are only part of the formal pipeline model. They may or may not correspond to real hardware signals.

The inner states and emitted signals of the units evolve in each cycle. The complexity of this state update varies from unit to unit. It can be as simple as a small table, mapping pending signals and inner state to a new state and signals to be emitted, e. g. for the IAG unit. It can be much more complicated, if multiple dependencies have to be considered, e. g. the instruction reconstruction and branch prediction in the IED stage. In this case, the evolution is formulated in pseudo code.

Full details on the model can be found in [8] and [12].

The output of the pipeline analysis is the number of cycles a basic block takes to execute, for each context. These results are then fed into the path analysis to obtain the WCET for the whole program.

The WCET analyzer gets as input:

• an executable (in ELF format). The code is generated with Diab Data C compiler from a restricted subset of ANSI-C (no dynamic data structures, no setjmp/longjmp),

™ WcetAl	? •
Analysis Options	
_ Executable	
Demo/data/mirmax.elf Browse	CFG
Start at:main	•
	1
About Analyze Visuali	ze <u>C</u> ancel

# Figure 4. WCET analyzer: Selection of the executable and start points.

- user annotations, giving the call targets for indirect function calls and upper bounds on the iteration counts of all loops,
- a description of the (external) memories and buses (i. e. a list of memory areas with minimal and maximal access times), and
- a task (identified by a start address). A task denotes a sequentially executed piece of code, no threads, no parallelism, and no waiting for external events. This should not be confused with a task in a operating system which might include code for synchronization or communication.

The WCET analyzer computes an upper bound of the runtime of the task (assuming no interference from the outside). Effects of interrupts, IO and timer (co-) processors are not reflected in the predicted runtime and have to be considered separately (e. g. by a quantitative analysis).

In addition to the raw information about the WCET, several aspects can be visualized by the **aiSee** tool [9] to view detailed information delivered by the analysis (see Figures 5-11).

#### 5 Conclusion

We presented a tool for the determining of the worst-case execution time for the Motorola ColdFire 5307. In con-



Figure 10. The timely development of pipeline states at a point in a program. The roots of the trees correspond to the 3 incoming pipeline states at this program point. Each layer in the trees corresponds to one CPU-cycle. Branches in the trees corresponds to unknown options, e.g. a cache hit and a cache miss at a memory access.



Figure 11. Detailed view pipeline states.



Worst Case Execution Time: 886

Figure 5. WCET analyzers results: Graphical representation of the call graph. The calls (edges) that contribute to the worst-case runtime are marked by the color red. The computed WCET is given in CPU cycles.

trast to other approaches (e. g. [2, 6, 7, 10, 13, 14, 15]) our tool takes into account the combinenation of all the different hardware characteristics while still obtaining tight upper bounds for the WCET of a given program in reasonable time. The tool was developed in the DAEDALUS project<sup>3</sup> and has been applied to a real-life benchmark containing realistically sized code modules. As recent trends, e. g., in automotive industries (X-by-wire, time-triggered protocols) require the knowledge of the WCET of tasks, such tools are of high importance.

We are currently finishing the pipeline analysis tool for the PowerPC 755, a processor that features out-of-order execution, speculation and superscalarity. Also, our technique opens the perspective to a generative approach, where analyses are generated from the specification of a model.



Figure 6. Basic block graph with runtime information for each block.



Figure 7. Unfolded basic block. For each instruction, the set of all possible pipeline states can be shown on demand.

<sup>&</sup>lt;sup>3</sup>The DAEDALUS project aims at introducing static program analysis methods into the airplane software validation process. Partners include Airbus France, AbsInt Angewandte Informatik GmbH, PolySpace and various academic partners, e. g. ENS, Saarland University, etc.







Figure 9. The worst case path trough the basic block graph of a routine is marked by the color red.

#### 6 Contact

For more information please contact:

AbsInt Angewandte Informatik GmbH Stuhlsatzenhausweg 69 D-66123 Saarbruecken Germany Phone: +49 681 8318317 Fax: +49 681 8318320 http://www.AbsInt.com info@AbsInt.com

#### References

- M. Alt and F. Martin. Generation of Efficient Interprocedural Analyzers with PAG. In *Proceedings of SAS'95, Static Analysis Symposium*, volume 983 of *Lecture Notes in Computer Science*. Springer-Verlag, 1995.
- [2] A. Colin and I. Puaut. Worst Case Execution Time Analysis for a Processor with Branch Prediction. *Real-Time Systems, Special issue on worst-case execution time analysis*, 18(2):249–274, 2000.
- [3] P. Cousot and R. Cousot. Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints. In *Fourth ACM Symposium on Principles of Programming Languages*, pages 238– 252, Los Angeles, California, 1977. ACM Press, New York.
- [4] C. Ferdinand. Cache Behavior Prediction for Real-Time Systems. PhD thesis, Universität des Saarlandes, 1997.
- [5] C. Ferdinand, R. Heckmann, M. Langenbach, F. Martin, M. Schmidt, H. Theiling, S. Thesing, and R. Wilhelm. Reliable and precise WCET determination for a real-life processor. In *Proceedings of EMSOFT 2001, First Workshop* on *Embedded Software*, volume 2211 of *Lecture Notes in Computer Science*, 2001.
- [6] C. A. Healy, R. D. Arnold, F. Mueller, D. B. Whalley, and M. G. Harmon. Bounding Pipeline and Instruction Cache Performance. *IEEE Transactions on Computers*, 48(1):53– 70, January 1999.
- [7] C. A. Healy, D. B. Whalley, and M. G. Harmon. Integrating the Timing Analysis of Pipelining and Instruction Caching. In *Proceedings of the IEEE Real-Time Systems Symposium*, pages 288–297, Dec. 1995.
- [8] R. Heckmann and S. Thesing. Cache and Pipeline Analysis for the ColdFire 5307. Technical report, Universität des Saarlandes, 2001.
- [9] http://www.aisee.com. aiSee Home Page.
- [10] Y. Hur, Y. H. Bae, S.-S. Lim, S.-K. Kim, B.-D. Rhee, S. L. Min, C. Y. Park, M. Lee, H. Shin, and C. S. Kim. Worst Case Timing Analysis of RISC Processors: R3000/R3100 Case Study. In *Proceedings of the IEEE Real-Time Systems Symposium*, pages 308–319, Dec. 1995.
- [11] M. Inc. MCF5307 ColdFire Integrated Microprocessor User's Manual. Motorola Inc., Aug. 2000. MCF5307UM/D, Rev. 2.0.

- [12] M. Langenbach, S. Thesing, and R. Heckmann. Pipeline Modeling for Timing Analysis. *Proceedings of the 9th In*ternational Static Analysis Symposium, 2002.
- [13] S.-S. Lim, Y. H. Bae, G. T. Jang, B.-D. Rhee, S. L. Min, C. Y. Park, H. Shin, K. Park, and C. S. Kim. An Accurate Worst Case Timing Analysis Technique for RISC Processors. *IEEE Transactions on Software Engineering*, 21(7), July 1995.
- [14] T. Lundqvist and P. Stenström. An Integrated Path and Timing Analysis Method based on Cycle-Level Symbolic Execution. *Real-Time Systems Journal*, 17(2/3):183–207, November 1999.
- [15] K. Narasimhan and K. Nilsen. Portable Execution Time Analysis for RISC Processors. 1994.
- [16] J. Schneider and C. Ferdinand. Pipeline Behavior Prediction for Superscalar Processors. Technical report, Universität des Saarlandes, May 1999.
- [17] J. Schneider and C. Ferdinand. Pipeline Behavior Prediction for Superscalar Processors by Abstract Interpretation. In Proceedings of the ACM SIGPLAN Workshop on Languages, Compilers and Tools for Embedded Systems, volume 34, pages 35–44, May 1999.
- [18] H. Theiling. Extracting Safe and Precise Control Flow from Binaries. In Proceedings of the 7th Conference on Real-Time Computing Systems and Applications, Cheju Island, South Korea, 2000.
- [19] H. Theiling. Generating Decision Trees for Decoding Binaries. In Proceedings of the ACM SIGPLAN Workshop on Language, Compiler and Tools for Embedded Systems, Snowbird, Utah, USA, June 2001.
- [20] H. Theiling. ILP-based Interprocedural Path Analysis. In Proceedings of the Workshop on Embedded Software, Grenoble, France, October 2002.
- [21] H. Theiling and C. Ferdinand. Combining Abstract Interpretation and ILP for Microarchitecture Modelling and Program Path Analysis. In *Proceedings of the 19th IEEE Real-Time Systems Symposium (RTSS)*, Madrid, Spain, 1998.

# The European Space Agency's involvement and interest is WCET and scheduling analysis

### **Extended Abstract**

Morten Rytter Nielsen, ESA (<u>morten.nielsen@esa.int</u>) Eric Conquet, ESA (<u>eric.conquet@esa.int</u>) Jean-Loup Terraillon (jean-loup.terraillon@esa.int)

### Abstract

We consider the use of scheduling analysis as not being a standalone exercise but a system-level activity, congruent with the conscious decision for a 'correctness by construction' development model. We describe how ESA have incorporated the ideas of scheduling analysis into our required standard practices; how we have ensured that the enabling technology is available; and where we see the future of WCET technology and scheduling analysis.

### The development approach

The traditional development model, which is used in software space projects under the responsibility of ESA, follows the classical waterfall V-model [ESA-PSS-05]. In this development model the User Requirements are ESA's requirements towards industry and the Software Requirements (or Technical Requirements) are industry's refinement of these. These Software Requirements are then followed by Architectural Design, Detailed Design and coding. On the ascending part of the V-model Unit Tests (verifies Detailed Design), Integration Tests (verifies Architectural Design), System Tests (verifies Software Requirements Definition) and Acceptance Tests (verifies User Requirements) are performed. The testing effort is usually 50-60 percent of the total development effort. Each phase of the development model is finalized with reviews and acceptance together with associated payments.

### Historic Space Systems

The V-model has proven its value through many years and projects. Traditionally onboard software-systems have been quite simple and with well separated functional blocks. The utilized software technology centered on fixed cyclic schedulers and dedicated proprietary kernels and very often the I/O mechanism was polling or well characterized interrupts. The required method in the ESA standards for controlling the performance behavior was limited to requirements for CPU utilization at the different stages of development (projects typically used 50% at Architectural Design, 60% at Detailed Design and 70% on final acceptance of the software code). Real-time requirements in the form of reactivity/responsiveness and jitter where either non-existent or at best occasional. The CPU utilization was typically acquired by estimation and later by measurement performed on the final code.

### **Current Space Systems**

The new generation of onboard space systems is significantly richer in functionality and complexity, with much more interaction between functional blocks than traditional onboard systems. Among other reasons, this trend originates from:

- Added throughput (dedicated services)
- More intelligent Autonomy and Failure, Detection, Isolation and Recovery (FDIR) functionality
- Intelligent instruments that sporadically interrupt the main computer
- Added capability of the onboard system in general

Many real-time requirements are now part of the requirement baseline to ensure reactivity and enable different units of the satellite to be developed to lesser tolerances.

Several new problems have surfaced in the new generation of onboard systems [ESA STR-260]. Many of these problems occur in the real-time behavior area. Since CPU utilization is not a sufficient way to ensure real-time behavior, the development approach has to be adapted. ESA have thus sponsored and funded a number of initiatives and supported (and still supports) the introduction of scheduling analysis in the ways outlined in the following paragraphs. For us it is clear that the use of scheduling analysis have major repercussions on the implementing technology as well as on the process standards and associated development approach. This altogether raises a clear demand for better tools support, not limited to the extraction of WCET and the scheduling analysis but also extending to the specification of the real-time attributes and properties of the system.

### Standards

The new European generation of space standards, the ECSS standards, allows more flexible development approaches to be used (e.g. spiral models and rapid prototyping) [ECSS-E40B-July2000 and ECSS-E40B-Feb2002]. However, they also require that the used computational model of the system be identified. This explicitly includes the component types (e.g. active-periodic, active-sporadic, protected, passive, actors and process), the assumed scheduling type and model (e.g. fixed priority or dynamic priority) and the accompanying analytical model under which the model is executed (e.g. Rate-Monotonic Scheduling).

This evolution shows that the previously informally used CPU utilization is now being replaced by much more stringent requirements on the chosen architecture and the rationale behind this choice.

Projects may decide to waive requirements in the standards if this implies too much effort. Thus the enabling technology is very important to lower the entrance to applying scheduling analysis.

### **Enabling technology**

#### **Specification and Design level**

In order to be able to really harvest the benefits of the scheduling technology early in the development process, ESA saw the need to accommodate the computational model already at design level. The result of this effort is the HOOD derived HRT-HOOD method [HRT-HOOD]. Currently, TNI (France) and Intecs Sistemi (Italy) have commercial tools supporting this specification and design method.

#### Implementation technology

ESA have supported the Ada Ravenscar definition from a user perspective. Furthermore we have funded the development of the GNAT/ORK kernel and compilation system and the port of Aonix Raven to the space processor ERC32. Also CNS have an Ada Ravenscar system for the ERC32. Ravenscar compilation systems are now used for Beagle2 and GOCE.

### WCET extraction

In some projects the extraction of the WCET profile have been done by hand. However, for scheduling analysis to be used widely and systematically in the space domain, we believe that tools supporting this process are needed. Various ways of acquiring the WCET have been tried, including:

- Instrumentation: Logic analyser (Tektronik) and embedded instrument code (Aonix and VxWorks/Tornado) plus user developed instrument code
- Source level analysis with the support of the compiler: a prototype based on the Adaworld compiler have been developed by Aonix
- Static Analysis on image code: Bound-T from SSF (Finland) have been developed for both the DSP 21020 and the ERC32

### Scheduling analysis

Tools to help apply different scheduling analysis techniques have been developed and are now available from Spacebell (Belgium). These tools assist in margin analysis and enables persons less fluent in the logic behind the analysis to interpret and evaluate the results.

### **Test cases**

A standardization of core onboard services has taken place in the form of the Packet Utilization Standard [ESA-PSS-05]. OBOSS [OBOSS] is a reference implementation of selected services, which have been used as a guinea pig for scheduling analysis and the Ada Ravenscar profile. Furthermore, the development approach using scheduling analysis and thereby moving the verification of real-time properties from the typical integration testing phase to the specification and design phase have been applied with great success on the European Robotic Arm (ERA) which is a safety critical module to be used on the International Space Station.

### Future of Space Systems

The new draft ECSS standards for onboard space engineering require that scheduling analysis must be performed. Several proposals for new onboard systems are baselining Ada Ravenscar as the implementation technology and the awareness of scheduling analysis is increasing. Together with standards that require a strong development baseline and a consolidation of the tools assisting in the scheduling analysis in all relevant phases of the development process, the entry barrier for the application of this development approach will be continuously lowered.

ESA continues to fund and promote the development of the enabling technology and the support for the development approach referenced in this paper. The near future evolution is the new space processor LEON, which like the current ERC32 has a Sparc instruction set. The transition to LEON, which has cache, raises new challenges that will require 'expert support' in addressing.

The movement and activities described in this paper has been triggered by problems encountered in space projects using the current development approach. These activities focus on a single computation platform with embedded software. As space onboard systems are moving from synchronous to asynchronous behavior, the need to extend the scheduling analysis to system level is surfacing. ESA is participating in organizations supporting the AADL (Avionic Architecture Description Language) standard. The aim of this work is to define a common language for the design and verification of complex avionic systems. We expect from such a standardization effort the emergence of an open framework that can incorporate various design languages and verification tools able to trap performance and behavioral issues in early design phases.

### Conclusion

We have in this extended abstract explained the context and the support of the WCET and scheduling analysis in ESA and the problems that we have encountered which let to this. In the full paper we will include experiences of the different areas outlined above and expand on the future as ESA sees it. This will include specific activities started or foreseen to be started in the area of distributed scheduling analysis.

### **References:**

**[ESA-PSS-05]** ESA PSS-05-0 Issue 2, February 1991: ESA Software Engineering Standards

**[ESA-PSS-07-101]** ESA PSS-07-101 issue 1, May 1994: Packet Utilisation Standard **[ECSS-E40B-July2000]** ECSS-E-40B Draft 1, 28 July 2000: Space Engineering. Software

**[ECSS-E40B-Feb2002]** ECSS-E40B Draft 1, 15 February 2002: Space Engineering. Software

**[HRT-HOOD]** Burns, A. and Wellings, A.: HRT-HOOD: A Structured Design Method for Hard Real-Time Ada Systems, Elsevier, 1995.

**[OBOSS]** OBOSS home page: http://spd-web.terma.com/Projects/OBOSS/Home\_Page **[ESA STR-260]** Vardanega, V.: Development of On-Board Embedded Real-Time Systems, ESA STR-260 October 1999

# Session IV: Low Level Analysis

#### Stefan Petters

### Presentation

This session was focussed on analysis issues of high performance processors. Isabelle Puaut started with a discussion of the advantages of locking data and code in the cache compared to using cache modelling techniques in her paper "*Cache Modelling vs Static Cache Locking for Schedulability Analysis in Multitasking Real-Time Systems*". The main results of this paper are that substantial gain in analysis complexity and time can be obtained and the WCET bound can be reduced considerably.

"A Framework to Model Branch Prediction for WCET Analysis" was presented by Tulika Nitra. This model in this paper focusses on the direct effects, i.e. the modelling of correct and incorrect predictions. Secondary effects like, for example, the preemption of a cache line due to a missprediction is left as an open issue for future work.

The paper "Difficulties in computing the WCET for Processors with Speculative Execution" of Pascal Sainrat describes the problems arising with speculative execution. Special focus is set on the fact that speculative execution dissolves the strict seperation of high level path analysis and low-level timing analysis utilised in static WCET analysis. Finally hints by which means this problems can be avoided are provided.

### 1 Discussion

As most embedded processors support cache locking the discussion on this issue circled around to what extend cache locking slows down or even speeds up the average execution of a task system. In small and highly critical systems the memory accesses and therefore the potential cache hits misses may be prefectly known. It was generally aggreed that the benefit of caches outweights the problems of it by far and therefor the caching processors should receive special attention in the research.

The branch prediction and corresponding speculative execution received a lot of attention in the discussion. The miss prediction ratio is usually less then 10 % and therefore the question arose whether the odd effects of a missprediction as, for example, cache lines beeing displaced by a missprediction are really worth the effort of detailed analysis or whether a more global view on things could improve the results.

An overall agreements was reached on the fact that the analysis of low-level effects for WCET analysis could be used to identify the hot spots of the code and therefore optimize the code to accelerate average execution times and make it better predictable.

### Cache analysis vs static cache locking for schedulability analysis in multitasking real-time systems

Isabelle Puaut INSA/IRISA, Campus de Beaulieu, 35042 Rennes Cédex, FRANCE e-mail: puaut@irisa.fr

#### Abstract

Cache memories have been extensively used to bridge the gap between high speed processors and relatively slow main memories. However, they are source of predictability problems and need special attention to be used in hard real-time systems. A lot of progress has been achieved in the last 10 years to model caches, in order to determine safe and precise bounds on (i) tasks WCETs in the presence of caches ; (ii) cache-related preemption delays. An alternative approach to cope with caches in real-time systems is to statically lock their contents so as to make memory access times and cache-related preemption times entirely predictable. This paper describes work in progress aiming at evaluating qualitatively and quantitatively the pros and cons of both classes of methods.

#### 1 Caches and real-time systems

Extensive studies have been performed on schedulability analysis to guarantee timing constraints in hard real-time systems. Schedulability analysis methods assume that task worst-case execution times (WCETs) are known. While many schedulability analysis methods consider that the cost of task preemption is zero to simplify the analysis, some methods account for task preemption costs (e.g. manipulation of task queues, cache-related preemption delays).

Caches are small and fast buffer memories used to speed up the memory accesses. They contain memory blocks that are likely to be accessed by the CPU in the near future. Although the caches are a very effective means of speeding up the memory accesses in the average case, they are a source of predictability problems, due to intra-task and inter-task interferences:

• *Intra-task* interferences occur when a task overrides its own blocks in the cache due to conflicts.

• *Inter-task* interferences arise in multitasking systems due to preemptions. The inter-task interferences imply a so-called *cache-related preemption delay* to reload the cache after a task is preempted.

Caches raise predictability issues in hard real-time systems because they are designed to speed up the system *average case* performance rather than the system *worst-case* performance which is of prime importance in hard real-time systems. As a consequence, the designers of hard real-time systems may choose not to use cache memories at all, or may choose to use on-chip static RAM – scratchpad memories – instead of caches [2]. The simple approach consisting in assuming that every access to memory results in a cache miss causes the tasks WCETs to be largely overestimated, which may cause the schedulability analysis to fail while the system may actually be feasible. The main issue is then to estimate tasks WCETs and cache-related preemption delays in a safe but not overly pessimistic manner.

Two classes of approaches, described hereafter, can be used to deal with caches in real-time systems.

**Cache analysis methods.** A first class of approaches to deal with caches in hard real-time systems is to use them without any restriction, and resort to *static analysis* techniques to predict their worst-case impact on the system schedulability.

At the intra-task level, static WCET analysis techniques have been extended to predict the impact of caching on the WCETs of the tasks. They achieve a classification of the memory accesses regarding the instruction or data caches (e.g. *hit* when it can be proved that the access always results in a cache hit, *miss* otherwise). Techniques to predict the worst-case task behavior regarding the instruction cache can use data-flow analysis on each task control flow graph [12], abstract interpretation [1], integer linear programming techniques [10], or symbolic execution [11]. At the inter-task level, work has been undertaken to obtain safe and precise estimates of the cache-related preemption delay [9]. In [9], at every possible preemption point, the blocks that will be used by each task after that point are determined by static analysis, thus avoiding considering that the whole memory accessed by the task has to be reloaded in the cache after a preemption.

**Cache partitioning and cache locking.** A second class of approaches to deal with caches in real-time systems is to use them in a restricted or customized manner, so as to adapt them to the needs of real-time systems and schedulability analysis.

Cache partitioning techniques [8, 5, 14] assign reserved portions of the cache (partitions) to certain tasks in order to guarantee that their most recently used code or data will remain in the cache while the processor executes other tasks. The dynamic behavior of the cache is kept within partitions. These techniques eliminate the inter-task interferences, but need extra-support to tackle intra-task interference (e.g. static cache analysis) and reduce the amount of cache memory available for each task.

Another way to deal with caches in real-time systems is to use *cache locking techniques*, which load the cache contents with some values and lock it to ensure that the contents will remain unchanged [6]. This ability to lock cache contents is available on several commercial processors. The cache contents can be loaded and locked at system start for the whole system lifetime (*static cache locking*), or changed during the system execution, like for instance when a task is preempted by another one (*dynamic cache locking*). The key property of cache locking is that the time required to access the memory is *predictable*.

Schedulability analysis for systems with caches. Some schedulability analysis methods have been extended to cope with cache-related preemption delays. They add the parameter  $\gamma_i$ , upper bound on the cache-related preemption delay, to the formulas in charge of verifying the system feasibility.

In [3], Rate Monotonic Analysis (RMA) is extended to cope with cache-related preemption delays. The utilization of a set of periodic tasks that takes the cache-related preemption delays into account is introduced (see equation 1 below).

$$U = \sum_{i=1}^{n} \frac{C_i + \gamma_i}{P_i} \tag{1}$$

In the equation, n is the number of tasks.  $C_i$  and  $P_i$  are the WCET and period of task number i. For static priority

systems with priorities assigned along the rate monotonic policy, a sufficient condition given in equation 2 below can then be used to verify the system schedulability [3].

$$U \le n(2^{\frac{1}{n}} - 1) \tag{2}$$

Response Time Analysis (RTA) has been extended by Busquets-Mataix et al [4] to take cache-related preemption delays into account, leading to the exact schedulability condition named CRTA. The principle of CRTA, for a task  $T_i$ , is to consider the interferences produced by the execution of the higher priority tasks on an increasing time window  $w_i^n$ . The response time  $R_i$  of task  $T_i$  is the fixed point of the sequence given in equation 3 below, with  $\gamma_i$  the cacherelated preemption delay and hp(i) the set of tasks that have a higher priority than  $T_i$ .

$$w_i^0 = C_i$$
  

$$w_i^{n+1} = C_i + \sum_{j \in hp(i)} \left\lceil \frac{w_i^n}{P_j} \right\rceil (C_j + \gamma_j) \to R_i \quad (3)$$

These series converge when  $\sum_{j \in hp(i) \cup \{i\}} \frac{C_j}{P_j} \leq 1$ . The response time  $R_i$  of task  $T_i$  can then be compared against its deadline to determine the schedulability of  $T_i$ .

#### 2 Cache analysis vs static cache locking

In the following, we give some elements that allow to choose between using statically locked caches or using the dynamic features of the caches together with static cache analysis techniques to bound accurately tasks WCETs and cache-related preemption delays. A static cache locking strategy with a frozen cache contents for *all tasks* is considered hereafter.

#### 2.1 Qualitative comparison

Static cache locking is attractive from several point of views. First of all, it improves the system performance compared to a system that does not use caches, with respect to both average and worst-case system performance.

In addition, with static cache locking, the time required to perform a memory access is *predictable* (it is either a *hit* or a *miss* depending on whether the value is locked in the cache or not). While WCET analysis is still required, it alleviates the need for using complex cache analysis techniques for computing WCETs and cache-related preemption delays, and results in more simple WCET analysis tools. In particular, it eliminates the issue of integrating cache analysis techniques with the analysis techniques for the other architectural features (pipelines, branch prediction, etc).

Static cache locking can also be used when no cache analysis method can apply, due for instance to nondeterministic or poorly documented cache replacement strategies (e.g. pseudo-random replacement policies).

Another important benefit of static cache locking is that the technique addresses both intra-task and inter-task interferences, which is unique among the cache management techniques presented above. Concerning inter-task interferences, since in static cache locking schemes the cache blocks are statically partitioned among the tasks, the cacherelated preemption delay is null, or is constant and equal to the time required to reload the processor prefetch buffer if the processor is equipped with such an architectural feature. This low cache-related preemption delay is particularly important for large caches (see section 2.2).

Finally, implementing cache locking turns out to be a light task once the contents of the locked cache are selected. No modification of the compilation process is required to implement static cache locking. In particular, the addresses of values (instructions/data structures) need not be modified, contrary to schemes that use static on-chip RAM to speed up memory accesses. To be implemented, static cache locking only requires to execute a small routine at the system start-up to load the contents of the cache with the selected values and lock the cache so that its contents remain unchanged during the whole system execution.

However, statically locking the contents of caches reduces the amount of cache memory available for each task. In addition, it raises the issue of selecting the cache contents. Since we are interested in hard real-time systems, the main objective of the cache selection algorithm is to improve the worst-case system behavior according to some of the metrics used by schedulability analysis methods, such as CPU utilization or interferences between tasks. The main issue is then to avoid performing an exhaustive search of all possible cache contents, which would require an untractable computation cost. For instance, for a direct-mapped cache, if up to 4 program lines can me mapped onto a given cache block, checking the feasibility of the system with all possible cache contents would require  $4^B$  feasibility tests, with B the number of cache blocks. This complexity led [6] to select a genetic algorithm for the selection of the cache contents and [13] to base the selection of cache contents on actual traces of the system execution.

Another potential benefit of static cache locking, although not proved yet by any study, is that it can easily apply to data caches, or to unified caches or to multi-level caches.

#### 2.2 Quantitative comparison

Since the primary focus in hard real-time systems is to prove that all deadlines are met, the key performance metric to be considered when comparing cache management schemes is the *worst-case* performance of the system. In this section, we compare the worst-case performance of a small task set made of periodic tasks using a state of the art *cache analysis technique* with its worst-case performance obtained using *static cache locking*. In this performance evaluation, we focus on *instruction caches* only.

#### 2.2.1 Experimental setup

**Target architecture and simulator.** The target architecture considered in the experiments is the simplified MIPS processor used by the Nachos operating system. Nachos<sup>1</sup> is a simple operating system for an emulated MIPS CPU, designed for teaching purposes. Nachos has been extended for the purpose of this study with an instruction cache with blocks of 16 bytes, and with a prefetch buffer of 16 bytes (4 MIPS instructions) to speed up sequential access to instructions. The instruction cache can be parametrized so that the cache size ranges from 512 bytes to 16Kbytes, and so that the associativity degree W ranges from 1 (direct-mapped cache) to 32 (set associative cache). The timing model considered for the processor is very simple: an instruction is assumed to execute in  $t_{hit} = 1$  clock cycle in the case of a cache hit, and  $t_{miss} = 10$  clock cycles otherwise.

**Static cache locking.** The algorithm designed to select the contents of the locked cache (see [13] for more details) applies to task sets made of *periodic* tasks. It aims at optimizing the task set schedulability, by minimizing the CPU utilization (as defined in equation 1,  $U = \sum_{i=1}^{n} \frac{C_i + \gamma_i}{P_i}$ ).

As the worst-case execution path required to compute the WCET  $C_i$  is not known unless the contents of the locked cache are known, the minimization of the utilization U is achieved thanks to the knowledge of an *actual* execution path for the tasks, which is not necessarily the *worst-case* path. This execution path is obtained through a simulation of the execution of the tasks using Nachos.

On the considered architecture, when static cache locking is used, the cache-related preemption delay  $\gamma_i$  is constant and equal to the delay required to refill the processor prefetch buffer ( $\gamma_i = t_{miss}$ ).

<sup>&</sup>lt;sup>1</sup>http://www.cs.washington.edu/homes/tom/nachos/

	Formulas for WCET computation: $W(S)$					
$S = S_1;; S_n$	$WCET(S) = WCET(S_i) + \dots + WCET(S_n)$					
$S = \text{if } (tst) \text{ then } S_1 \text{ else } S_2$	$WCET(S) = WCET(Test) + \max(WCET(S_1), WCET(S_2))$					
$S = \log(tst)S_1$	$WCET(S) = maxiter * (WCET(Test) + WCET(S_1)) + WCET(Test)$					
	where <i>maxiter</i> is the loop maximum number of iterations.					

Table 1. Simplified equations for tree-based WCET computation

Task name	Description	Code size	WCET-	Period
		(Bytes)	miss	
qurt	Computation of roots of quadratic equations	1824	21474	59697
minver	Matrix inversion	4320	36701	70098
jfdctint	JPEG integer implementation of the forward DCT	3440	29324	127559
fft1	FFT (Fast fourier transform) Cooly-Turkey algorithm	3620	115152	601093

Table 2. Task set characteristics

WCET analysis and static cache analysis. Worst-case execution times (WCET) are obtained using the Heptane tree-based WCET analysis tool [7]. Heptane computes WCETs through a bottom-up traversal of the syntax tree of the analyzed programs. Table 1 gives a flavor of the formulas used in Heptane to compute WCETs (these formulas are overly simplified since they assume constant and context-independent execution times for the leaves of the syntax tree).

Heptane includes hardware modeling capabilities so as to estimate safely but precisely the WCETs on architectures with instruction caches, pipelines and simple branch predictors. Here, Heptane's pipeline and branch prediction modeling modules have been switched off since our focus is on instruction caches only.

The technique used in Heptane to estimate the instruction cache behavior is based on F. Mueller's so-called static cache simulation [12]. The cache analysis technique computes abstract cache states (representation of all possible cache contents considering all possible execution paths in the program), computed using data-flow analysis on the program control flow graph. Abstract cache states are then used to classify the instructions according to their worst-case behavior regarding the instruction cache (e.g. *hit* when it is certain that there are no conflict for a cache block, miss when no better classification can be found, and intermediate categories when there are conflicts inside loops). In the case of set-associative caches, when a given instruction appears in different blocks of the same set of an abstract cache state, different categories can be obtained for the different instances of the instructions. In this situation, in order for the analysis to be safe, the instruction is classified with the most pessimistic category.

Heptane has been modified so as to offer the possibility to replace the cache analysis module by a module that takes into account the presence of locked caches. This new module classifies instructions into two categories: *miss* and *hit*. An instruction is classified as a *hit* if it is locked in the instruction cache, and is classified as a *miss* otherwise.

When Heptane is configured to be used with its cache analysis capabilities, no attempt is made to bound the cacherelated preemption delay  $\gamma_i$  precisely (it is assumed that all program lines of a given task have to be reloaded after a preemption, with a maximum of N reloads where N is the number of cache lines).

**Task set.** The worst-case performance evaluation has been achieved on a small real (non synthetic) task set, whose characteristics are given in table 2.

The table gives for every task: its name, a short description, the size of its code in bytes, its WCET assuming that all instruction fetches cause a cache miss, and its period. The delays in the table are expressed in number of processor cycles. The periods of tasks have been selected so that the CPU utilization  $(\sum_{i=1}^{n} \frac{C_i + \gamma_i}{P_i})$  equals 1.3 (i.e. the task set is not feasible if no instruction cache is used).

#### 2.2.2 Worst-case performance analysis

The worst-case system performance of the considered task set is given in Table 3. Each cell indicates whether the task set is feasible or not according to CRTA (equation 3, detailed in [4]). A '+' sign means that the task set if feasible, whereas a '-' sign means that it is not. The CPU utilization

Asso	Size	512B	1KB	2КВ	4KB	8KB	16KB
1	Locking	-1.188	-1.073	+0.936	+0.774	+0.740	+0.740
	Analysis	+0.785	+0.605	+0.550	+0.565	+0.567	+0.567
2	Locking	-1.174	-1.039	+0.881	+0.699	+0.551	+0.516
	Analysis	+0.833	+0.701	+0.588	+0.575	+0.567	+0.567
4	Locking	-1.169	-1.022	+0.844	+0.678	+0.506	+0.395
	Analysis	+0.935	+0.806	+0.687	+0.620	+0.577	+0.567
8	Locking	-1.177	-1.030	+0.823	+0.644	+0.495	+0.393
	Analysis	+0.943	+0.954	+0.813	+0.718	+0.623	+0.577
16	Locking	-1.177	-1.052	+0.825	+0.641	+0.489	+0.393
	Analysis	+0.921	-0.997	-1.025	+0.840	+0.721	+0.623
32	Locking	-1.176	-1.044	+0.817	+0.625	+0.487	+0.393
	Analysis	+0.904	-1.031	-1.113	-1.057	+0.842	+0.721

Table 3. Compared worst-case performance of static cache locking and static cache analysis

of the task set (as defined in equation 1) is also given in each cell. These two pieces of information are given for different cache sizes (Bytes), degrees of associativity, and this with static cache analysis (label *Analysis* in the table) and with static cache locking (label *Locking* in the table).



Figure 1. Worst-case CPU utilization

Figure 1 depicts the CPU utilization obtained on the task set from the contents of table 3. It compares the CPU utilization obtained when using cache locking and static cache analysis.

**Performance for small cache sizes.** It can be noted from the contents of figure 1 and table 3 that for small cache sizes, static cache analysis performs better than static cache locking, both from the standpoint of the task set feasibility and from the standpoint of CPU utilization. This comes from the fact that for small caches, there is a high degree of intra-task and inter-task conflicts for every cache block. By construction, the static cache analysis used in the WCET analyzer is not sensitive to inter-task interferences because tasks are considered separately. Thus, the WCET of tasks with cache analysis is lower than with cache locking. Moreover, the comparatively lower WCETs of cache analysis come with a higher cache-related preemption delay. However, the impact that this higher cache-related preemption delay has on the system schedulability and CPU utilization is negligible for small caches.

**Impact of the cache size.** It can be noted that for a given degree of associativity, the performance of both static cache locking and static cache analysis increases with the cache size, because of the decrease of the number of conflicts for cache blocks. However, the performance increase of static cache locking is higher than the one of static cache analysis when the cache size increases. For instance, for a 4-way associative cache, the task set exhibits better (i.e. lower) CPU utilization when using static cache locking than when using cache analysis for caches larger than 8KB, whereas static cache locking performs better for caches smaller than 8KB. This is because the cache-related preemption delay increases linearly with the cache size for the static cache analysis method, whereas it stays constant for the static cache locking method.

Impact of the degree of associativity. For a given cache size, the performance of static cache locking scales better than the one of static cache analysis with an increasing degree of associativity W. Indeed, static cache locking takes

benefit of the increasing degree of associativity to eliminate both intra-task and inter-task interference, which explains that the CPU utilization increases with W. In contrast, the static cache analysis method we have used does not scale well with W. This comes from the *pessimistic* way the instructions are classified (see § 2.2.1).

#### **3** Open issues

The key benefits of static cache locking is to make the time required to perform memory accesses predictable, and to be a unified technique to take into account both intratask and inter-task conflicts for cache blocks. This class of techniques alleviates the need for using complex static analysis techniques for computing WCETs and cache-related preemption delays. In addition, it can be applied in situations where static cache analysis cannot be used at all (*e.g.* when the instruction cache has a non deterministic or non documented cache replacement policy). While algorithms already exist for selecting the contents of statically locked caches [6, 13], we think that further work is required:

- to study their performance on larger real (non synthetic) benchmarks, in particular in task sets whose size is much larger than the cache size. For large programs, a possible direction is to explore more dynamic cache locking strategies (for instance, to select different contents of the locked cache changed at staticallydefined points in order to cope with the tasks dynamic behavior while staying predictable)
- to study the impact of statically locked caches on the system average case performance
- to study the applicability of static cache locking techniques to data/unified/multi-level caches
- to address implementation issues on actual embedded processors
- to compare the use of statically locked caches with the use of on-chip static RAMs (benefits wrt predictability, issues to be addressed)

#### Acknowledgements

Thanks to David Decotigny (IRISA) and Jörn Schneider (Saarland University) for pointing out errors in earlier drafts of this paper.

#### References

[1] M. Alt, C. Ferdinand, F. Martin, and R. Wilhelm. Cache behavior prediction by abstract interpretation. In SAS'96, Static Analysis Symposium, volume 1145 of Lecture Notes in Computer Science, pages 51–66. Springer, September 1996.

- [2] O. Avissar, R. Barua, and D. Stewart. Heterogeneous memory management for embedded systems. In Proc. of the International Conference on Compilers, Architectures and Synthesis for Embedded Systems, Atlanta, GA, USA, Nov. 2001.
- [3] S. Basumallick and K. Nilsen. Cache issues in real-time systems. In ACM SIGPLAN Workshop on Languages, Compilers, and Tools for Embedded Systems, June 1994.
- [4] J. V. Busquets-Mataix, J. J. Serrano, R. Ors, P. Gil, and A. Wellings. Adding instruction cache effect to schedulability analysis of preemptive real-time systems. In *Proceedings* of the 1996 Real-Time technology and Applications Symposium, pages 204–212. IEEE Computer Society Press, June 1996.
- [5] J. V. Busquets-Mataix and A. Wellings. Hybrid instruction cache partitioning for preemptive real-time systems. In *Proc. of the 9th Euromicro Workshop of Real-Time Systems*, pages 56–63, Toledo, Spain, June 1997.
- [6] M. Campoy, A. P. Ivars, and J. V. Busquets-Mataix. Static use of locking caches in multitask premptive real-time systems. In *IEEE/IEE Real-Time Embedded Systems Workshop* (Satellite of the IEEE Real-Time Systems Symposium), London, UK, Dec. 2001.
- [7] A. Colin and I. Puaut. A modular and retargetable framework for tree-based wcet analysis. In *Proc. of the 13th Euromicro Conference on Real-Time Systems*, pages 37–44, Delft, The Netherlands, June 2001.
- [8] D. B. Kirk. Smart (strategic memory allocation for realtime) cache design. In *Proceedings of the 10th IEEE Real-Time Systems Symposium (RTSS89)*, pages 229–237, Santa Monica, California, USA, Dec. 1989.
- [9] C. G. Lee, J. Hahn, Y. M. Seo, S. L. Min, R. Ha, S. Hong, C. Y. Park, M. Lee, and C. S. Kim. Analysis of cache-related preemption delay in fixed-priority preemptive scheduling. *IEEE Transactions on Computers*, 47(6), June 1998.
- [10] Y.-T. S. Li, S. Malik, and A. Wolfe. Cache modeling for real-time software: Beyond direct mapped instruction cache. In *Proceedings of the 17th IEEE Real-Time Systems Symposium (RTSS96)*, pages 254–263. IEEE, IEEE Computer Society Press, Dec. 1996.
- [11] T. Lundqvist and P. Stenstrom. An integrated path and timing analysis method based on cycle-level symbolic execution. *Real-Time Systems*, 17(2-3):183–207, Nov. 1999.
- [12] F. Mueller. Timing analysis for instruction caches. *Real-Time Systems*, 18(2):217–247, May 2000.
- [13] I. Puaut and D. Decotigny. Low-complexity algorithms for static cache locking in multitasking hard real-time systems. Submitted to publication - available on demand, May 2002.
- [14] J. E. Sasinowski and J. K. Strosnider. A dynamic programming algorithm for cache/memory partitioning for realtime systems. *IEEE Transactions on Computers*, 42(8):997– 1001, Aug. 1993.

### A Framework to Model Branch Prediction for WCET Analysis

Tulika Mitra Department of Computer Science School of Computing National University of Singapore Singapore 117543 tulika@comp.nus.edu.sg

In this paper, we present a framework to model branch prediction for Worst Case Execution Time (WCET) analysis. Our micro-architectural modeling is completely generic, and parameterizable w.r.t. the currently used branch prediction schemes. It automatically derives linear constraints on the total misprediction count from the control flow graph of the program. These constraints can be solved by any integer linear programming (ILP) solver to estimate the WCET.

Current generation processors perform control flow speculation through branch prediction, which predicts the outcome of branch instructions. If the prediction is correct, then execution proceeds without any interruption. For incorrect prediction, the speculatively executed instructions are undone, incurring a branch misprediction penalty between 3-19 clock cycles. If branch prediction is not modeled, all the branches in the program must be conservatively assumed to be mispredicted for finding the WCET. This pessimism results in as much as 60 - 70% overestimation for some of the benchmarks in this paper, even assuming a 3 clock cycle branch misprediction penalty.

A classification of branch prediction schemes appears in Figure 1. Branch prediction can be *static* or *dynamic*. Static schemes associate a fixed prediction to each branch instruction via compile time analysis. Almost all modern processors, however, predict the branch outcome dynamically based on past execution history. Dynamic schemes are more accurate

Abhik Roychoudhury Department of Computer Science School of Computing National University of Singapore Singapore 117543 abhik@comp.nus.edu.sg



Figure 1: Classification of Branch Prediction Schemes. At each level, the more widely used category is underlined.

than static schemes, and in this work we study only dynamic branch prediction. The first dynamic technique proposed is called *local branch prediction* [4], where each branch is predicted based on its last few outcomes. This scheme uses a  $2^n$ -entry *branch prediction table* to store the past branch outcomes, which is indexed by the *n* lower order bits of the branch address. In the simplest case, each prediction table entry is 1-bit and stores the last outcome of the branch mapped to that entry. When a branch is encountered, the corresponding table entry is looked up and used as the prediction. When a branch is resolved, the corresponding table entry is updated with the outcome. A more accurate version of local scheme uses *k*-bit counter per table entry.

Most modern processors however use *global* branch prediction schemes [4] (also called correlation based schemes), which are more accurate. Examples of processors using global branch prediction include Intel Pentium Pro, AMD, Alpha as well as embedded processors IBM PowerPC 440GP and SB-1 MIPS 64. In these schemes, the prediction of the outcome of a branch I not only depends on I's recent outcomes, but also on the outcome of the other recently executed branches. Global schemes can exploit the fact that behavior of neighboring branches in a program are often correlated. Global schemes uses a single shift register, called branch history register (BHR) to record the outcomes of n most recent branches. As in local schemes, there is a global branch prediction *table* in which the predictions are stored. The various global schemes differ from each other (and from local schemes) in the way the prediction table is looked up when a branch is encountered.

Little work has been done to study the effects of branch prediction on WCET. Effects of static branch prediction have been investigated in [1, 3]. However, most current day processors (Intel Pentium, AMD, Alpha, SUN SPARC) implement dynamic branch prediction schemes, which are more difficult to model. To the best of our knowledge, [2] is the only other work on timing estimation under dynamic branch prediction. However, their technique only models the effects of local prediction schemes.

The starting point of our analysis is the control flow graph (CFG) of the program. Let  $v_i$  denote the number of times block i is executed, and let  $e_{i,j}$  denote the number of times control flows through the edge  $i \rightarrow j$ . As inflow equals outflow,  $v_i = \sum_{j \rightarrow i} e_{j,i} = \sum_{i \rightarrow j} e_{i,j}$ . We provide bounds on the maximum number of iterations for loops and maximum depth of recursive invocations for recursive procedures. These bounds can be user provided, or can be computed off-line for certain programs.

Let  $cost_i$  be the execution time of basic block iassuming perfect branch prediction. Given the program,  $cost_i$  is a fixed constant for each i. Then, the total execution time of the program is  $\sum_i (cost_i * v_i + penalty * m_i)$  where *penalty* is a constant denoting the penalty for a single branch misprediction;  $m_i$  is the number of times the branch in block i is mispredicted. By maximizing this objective function we can get WCET.

**Modeling Prediction Schemes** To determine the prediction of a block *i*, we first compute the index into the prediction table. We define  $v_i^{\pi}$  and  $m_i^{\pi}$ : the execution count and the misprediction count of block *i* when branch in *i* is executed with index =  $\pi$ . By definition:

$$m_i^{\pi} \le v_i^{\pi}$$
  $m_i = \sum_{\pi} m_i^{\pi}$   $v_i = \sum_{\pi} v_i^{\pi}$ 

The prediction schemes differ from each other primarily in how they index into the prediction table. To predict a branch I, the index computed can be a function of: (a) the past execution trace (history) and (b) address of the branch instruction I. In the GAgscheme, the index computed depends solely on the history and not on the branch instruction address. Other global prediction schemes (gshare, gselect) use both history and branch address, while local schemes use only the branch address.

Our modeling is independent of the definition of the prediction table index  $\pi$ . Hence it can apply to any branch prediction scheme that uses a single prediction table. To model the effect of different branch prediction schemes, we only alter the meaning of  $\pi$ , and show how  $\pi$  is updated with the control flow.

In the case of GAg, this index is the outcome of last k branches before block i is executed. These k outcomes are recorded in the Branch History Register (BHR). To model the change in history due to control flow, we use the left shift operator; thus  $left(\pi, 0)$ shifts pattern  $\pi$  to the left by one position and puts 0 as the rightmost bit. We define:

**Definition 1** Let  $i \to j$  be an edge in the control flow graph and let  $\pi$  be the BHR content at basic block *i*. The change in history pattern on executing  $i \to j$  is given by  $\Gamma(\pi, i \to j) = \pi$  if  $i \to j$  is an unconditional jump. If  $i \to j$  is a taken (non-taken) branch then  $\Gamma(\pi, i \to j)$  is  $left(\pi, 0)$  ( $left(\pi, 1)$ ).

In the popular gshare [4] scheme, the BHR is XORed with last n bits of the branch address to look up the prediction table. Usually, gshare results in

Pgm.	gshare		GAg		local	
	Mispred		Mispred		Mispred	
	Obs.	Est.	Obs.	Est.	Obs.	Est.
check	3	3	3	3	198	198
matsum	204	204	204	204	200	200
matmul	223	223	223	223	200	200
fdct	7	7	7	7	4	4
fft	3678	6165	3398	5175	4129	5154
isort	9687	9952	587	598	399	399
bsearch	9	9	9	10	6	7
eqntott	203	205	202	206	203	204

Table 1: Observed and estimated misprediction count with gshare, GAg, and local schemes.

a more uniform distribution of table indices compared to *GAg.* We define the index  $\pi$  as  $\pi$  =  $history_m \oplus address_n(I)$  where m, n are constants,  $n \ge m, \oplus$  is XOR,  $address_n(I)$  denotes the lower order n bits of I's address, and  $history_m$  denotes the most recent m branch outcomes (which are XOR-ed with higher-order m bits of  $address_n(I)$ ). And,

$$\Gamma_{gshare}(\pi, i \to j) = \Gamma(history_m, i \to j) \oplus address_n(I)$$

In local schemes, the index  $\pi$  for branch instruction I is the least significant n bits of I's address, denoted  $address_n(I)$  (n is a constant). Here  $\pi$  is independent of the past execution history of other branches. The update of  $\pi$  due to control flow is given by  $\Gamma_{local}(\pi, i \to j) = address_n(J)$ , where  $address_n(J)$  denotes the least significant n bits of the last instruction J in basic block j.

**Bounding Mispredictions** Given the definition of  $\pi$  and  $\Gamma$ , we derive inflow and outflow constraints on the flow of  $\pi$  through the control flow graph to derive upper bounds on  $v_i^{\pi}$ . To bound  $m_i^{\pi}$ , we note the following. Suppose there is a misprediction of the branch in block *i* with history  $\pi$ . This means that certain blocks (maybe *i* itself) were executed with history  $\pi$ , the outcome of these branches appear in the  $\pi$ th row of the prediction table, and the outcome of these branches *must have created* a prediction different from the current outcome of block *i*. To model mispredictions, we therefore capture repeated occurrence of a history  $\pi$  during program execution with differing outcomes; we provide constraints to bound such occurrences. Details of our modeling appear in [5] and are ommitted here for space considerations.

**Experimental Results** We selected eight different benchmarks for our experiments. We assumed zero cache misses and a perfect processor pipeline with no stalls except for penalty due to misprediction of conditional branches. These assumptions, although simplistic, allow us to separate out and measure the accuracy of our estimation technique. We assumed that the branch misprediction penalty is 3 clock cycles (as in the Intel Pentium processor). We used the SimpleScalar architectural simulation platform in the experiments. By changing SimpleScalar parameters, we could change the branch prediction scheme for the experiments.

To evaluate the accuracy of our branch prediction modeling, we present the experiments for three different branch prediction schemes: gshare, GAg and local. Since finding the worst case input of a benchmark (which produces the actual WCET) is a human guided and tedious process, we only measured the actual WCET assuming a 4-entry prediction table. The results appear in Table 1. In this table, we have shown only the observed and estimated misprediction counts to enable clear understanding of the accuracy of our technique (which models the effect of branch prediction). Even though not shown here due to space shortage, the estimation accuracy was independent of the prediction table size. Our estimation technique obtains a very tight bound on the WCET and misprediction count in all benchmarks except fft. The reason is that the number of iterations of the innermost loop of fft depends on the loop iterator variable value of the outer loops. This problem can be solved by providing expressions on the loop iteration counts instead of constants, as shown in [2].

Using CPLEX, a commercial ILP solver distributed by ILOG, on a Pentium IV 1.3 GHz processor with 1 GByte of main memory, our timing estimation technique requires less than 0.11 second for all the benchmarks with prediction table size varying 4–1024 entries.

One major concern with any ILP formulation of WCET is the scalability of the resulting solution. To check the scalability of our solution, we formulated the WCET problem for the popular gshare scheme with branch prediction table size varying from 4–1024 entries. Recall that in *gshare*, the branch instruction address is XOR-ed with the global branch history bits. In practice, gshare scheme uses smaller number of history bits than address bits, and XORs the history bits with the higher order address bits [4, 6]. The choice of the number of history bits in a processor depends on the expected workload. In our experiments, we used a maximum of 4 history bits as it produces the best overall branch prediction performance across all our benchmarks. As Figure 2 shows, the number of variables generated for the ILP problem initially increases and then decreases. With increasing number of history bits, number of possible patterns per branch increases. But with fixed history size and increasing prediction table size, the number of cases where two or more branches have the same pattern starts to decrease. This significantly reduces the number of ILP variables.

### References

 K. Chen, S. Malik, and D.I. August. Retargatable static software timing analysis. In *IEEE/ACM Intl. Symp. on System Synthesis (ISSS)*, 2001.



Figure 2: Change in ILP problem size with increase in number of entries in the branch prediction table for *gshare* scheme

- [2] A. Colin and I. Puaut. Worst case execution time analysis for a processor with branch prediction. *Journal of Real time Systems*, May 2000.
- [3] S-S. Lim, J.H. Han, J. Kim, and S.L. Min. A worst case timing analysis technique for inorder superscalar processors. Technical report, Seoul National University, 1998. Earlier version published in IEEE Real Time Systems Symposium(RTSS) 1998.
- [4] S. McFarling. Combining branch predictors. Technical report, DEC Western Research Laboratory, 1993.
- [5] T. Mitra and A. Roychoudhury. Effects of branch prediction on worst case execution time of programs. Technical Report 11-01, School of Computing, National University of Singapore, 2001.
- [6] S. Sechrest, C-C. Lee, and T. Mudge. Correlation and aliasing in dynamic branch predictors. In ACM International Symposium on Computer Architecture (ISCA), 1996.

### Difficulties in Computing the WCET for Processors with Speculative Execution

Christine Rochange and Pascal Sainrat Institut de Recherche en Informatique de Toulouse, France {rochange, sainrat}@irit.fr

#### Abstract

In real-time applications, the Worst-Case Execution Time often needs to be estimated to check that deadlines will be respected. With the trend of using up-to-date processors, WCET computation techniques continuously have to evolve in order to take into account the most recent hardware features. In this paper, we show that speculative execution ignoring can lead to underestimated execution times, and we explain why modelling it is not straightforward. We feel that pure static analysis might not allow safe WCET computation, due to the fact that speculative execution prevents the decoupling between the high-level (path) analysis and the low-level (timing) analysis.

#### 1. Introduction

For a large class of applications, embedded software has to satisfy hard real-time constraints. This requires to be able to tightly estimate the worst-case execution time (WCET) of programs.

WCET analysis has received much attention these ten last years. Dynamic methods involve measurements on real hardware or on cycle-level simulators. All the possible execution paths have to be explored in order to obtain the longest execution time. This poses two problems: (i) the number of possible paths is generally high and then the measurement time is prohibitive; (ii) for each path, the corresponding input data set has to be defined, which is usually difficult. In response to the drawbacks of dynamic methods, several static approaches have been proposed. They consist in three steps. First, the high-level analysis considers the program code in order to identify the possible execution paths, where a path is a list of basic blocks. Second, the low-level analysis estimates the execution time of each basic block. It is carried out in two phases: the global low-level analysis takes into account hardware components the behaviour of which depends on the global history of execution (e.g. cache memories); the local phase models components that only depends on the recent history (e.g. pipeline). Third, the execution times of paths are computed and the WCET is the longest one.

However, embedded systems tend to use modern processors featuring advanced architectural mechanisms that might be hard to model. Among these mechanisms, branch prediction, sometimes coupled with speculative execution, is implemented in most of the recent processors.

Estimating the WCET for processors with speculative execution does not present any special difficulty when it is based on dynamic measures: either the real target hardware is available (with speculative execution activated), or a cycle-level simulator is used, and speculative execution is not harder to model than other advanced features. However, current dynamic measurement methods often require to explore a too large number of execution paths and, for this reason, static analysis is generally preferred.

In this paper, we will show that estimating the WCET when a processor implements speculative execution is not straightforward. We suggest that usual static analysis techniques might not allow safe WCET computation, highlighting situations where they would lead to underestimation of the execution time.

Section 2 gives an overview of branch prediction and speculative execution techniques, and presents the work of Colin and Puaut [1] that takes branch prediction (but not speculative execution) into account within static WCET analysis. Section 3 shows why it is important to carefully model speculative execution to obtain a safe WCET. Section 4 discusses the difficulties of doing it within pure static WCET analysis, and section 5 concludes the paper.

# **2.** Branch prediction and speculative execution

#### 2.1 Overview

Modern processors are designed around longer and longer pipelines. Whenever a branch instruction is encountered in the instruction flow, the correct execution path is not known until the branch is executed. To avoid interrupting the instruction fetching, one of the two possible paths is speculatively selected by a branch predictor and instruction processing continues along this path. When the branch is resolved and if the speculative path is not the correct one, recovery actions are taken (e.g. the pipeline is flushed) and instruction processing restarts from the branch along the right path. Processing along a speculative path means fetching instructions from
the memory hierarchy, decoding and dispatching them to the reservation stations where they wait for their operands. For a processor that implements out-of-order execution, instructions belonging to the speculative path might also be executed before earlier instructions, and in particular before unresolved branches. This is what is called *speculative execution*. In that case, recovery from branch misprediction is a bit more complicated and generally requires mechanisms to restore the correct architectural state. Note that recovery is only required for components that must have a safe behaviour: the effects of a branch prediction error on other components, like cache memories or the branch predictor itself do not endanger correct functional results, they only might lower the system performance.

Many algorithms exist to predict the issue of branch instructions. The most recent ones include three kinds of structures:

- the *PHT* (*Pattern History Table*) is used to predict the direction of conditional branches: each of its entries reflects a recent history (often as a 2-bit saturating counter). The PHT is usually not tagged and it can be indexed by the instruction address (PC) alone or combined with a global or a local history recorded in one or several *BHR* (*Branch History Register*). Thus, several branches share the same counter and, on the contrary, the behaviour of a branch depends on several counters according to the history.
- the *BTB* (*Branch Target Buffer*) is used to predict the target address (except for subroutine returns)
- the *RAS* (*Return Address Stack*) is used to predict subroutine returns.

## 2.2 Computing the WCET for processors with branch prediction

As far as we know, branch prediction has been considered within WCET analysis only for in-order processors: this work has been presented by Colin and Puaut [1]. They consider the Intel Pentium, which features a simple branch predictor based on a single table referred to as BTB. The proposed method includes several stages.

First, the control-flow graph is analysed to build an *abstract state* of the BTB for each basic block: it indicates which instructions might be contained in each entry of the BTB before and after the execution of the basic block. The *input* abstract state of a basic block is computed from the *output* abstract states of the possible preceding basic blocks. Then, the abstract states are used to classify the branch instructions and to determine, for each of them, if it will be correctly predicted or not.

When ever this cannot be statically decided, the instruction is assumed to be mispredicted, which is supposed to be the worst case.

The WCET is then computed in two steps. First, a perfect branch predictor is assumed, and the WCET is estimated from the syntax tree and a set of formulas that express the maximum execution time of algorithmic structures. Then the timing effects of prediction errors are evaluated for a real branch predictor: a penalty delay is associated to each possibly mispredicted branch instruction. A second set of formulas is used to recursively build delay sets for each algorithmic structure of the syntax tree. The sum of these delays is then added to the WCET previously computed with perfect branch prediction.

#### 3. Possible effects of speculative execution

When a processor implements speculative execution, processing along the wrong path may have two kinds of effects on the system. In this section, we describe these effects and show why ignoring them can lead to underestimate the WCET.

#### 3.1 Dynamic instruction scheduling

Instructions of the wrong path occupy hardware resources, like functional units. Now, some flushing policies implemented for branch prediction error recovery do not immediately free the functional units. Thus, an instruction of the wrong path might continue its execution in a multi-cycle functional unit after the flushing of the pipeline (but its result will then be simply discarded). If the functional unit is not pipelined, the execution of later instructions belonging to the correct path might be delayed. Then the misprediction penalty would be longer than the strict recovery time.

Moreover, inserting wrong path instructions in the pipeline can modify the scheduling of previous instructions. For example, an instruction belonging to the wrong path that has its operands ready can be scheduled before a preceding instruction that is waiting for one of its operands. This might completely modify the overall scheduling, and then the execution time as mentioned in [3].

If speculative execution is not taken into account, the pipeline reservation tables produced by the local timing analysis might not be correct and the computed WCET could be underestimated.

#### 3.2 Memory contents

Processing along the wrong path can also change the content of memories. If instructions of the wrong path miss in the instruction cache, they are fetched from the upper level of the memory hierarchy. This can have a detrimental effect on the program execution time if instructions of the wrong path replace in the cache instructions belonging to the correct path: when the execution later restarts along the right path, those replaced instructions will miss in the cache, thus requiring longer fetch times. This detrimental effect is often referred to as *cache pollution*. Note that fetching instructions along the wrong path can also have a beneficial effect, as reported in [4]: some instructions of the wrong path can later be found on the correct path, and then processing along the wrong path acts as a prefetch mechanism.

The same effects can be observed on data accesses, provided that instructions of the wrong path are executed (not only fetched), which is only allowed in dynamicallyscheduled processors.

Processing along the wrong path may also have a beneficial or detrimental impact on the memories of the branch predictor (BTB, BHR, PHT and RAS) if they are updated speculatively in the earlier stages of the pipeline [2]. Only few parts are checkpointed for cost reasons (e.g. checkpointing the BTB is probably not affordable). If recovery is not implemented, the branch predictor tables might be polluted by the execution of the wrong path.

Now, let us assume that, ignoring speculative execution, the global low-level analysis is able to statically determine the real behaviour of all instruction and data cache accesses (hit or miss) and of all branches (well- or wrong- predicted). To understand why the possible pollution of memories due to wrong path execution should not be ignored, let us consider the following example:

This program may be compiled as:

```
T.0
      i = 0:
L1
      if i==10 then branch to L7
T.3
      s[i]=0
      i = 0
      if i==10 then branch to L6
T.4
L5
      s[i] = s[i] + t[i][j]
      i++
      branch to L4
L6
      m[i] = s[i] / 10
      i++
      branch to L1
L7
```

If we consider a branch prediction algorithm based on 2-bit saturating counters, initialised to "weakly-taken", the branch instructions of basic blocks L1 and L4 are mispredicted in the first iteration of loops i and j, and correctly predicted in the other iterations, while those of basic blocks L5 and L6 are always well predicted. As far as data accesses are concerned, the first reference to s[i] is determined to miss in the data cache while the other ones should hit.

Now, what does really happen if the processor implements speculative execution? At the first iteration of loop i, since the branch of basic block L1 is mispredicted, some instructions belonging to the wrong path are processed. In particular, access to m[i] might be executed. If m[i] happens to fall in the same cache line as s[i] then, when the branch is resolved and the execution restarts along the correct path, s[i] misses in the data cache, contrarily to the conclusion of the global low-level analysis. As a result, the actual execution time might be longer than the estimated WCET, which is not acceptable.

In the same manner, the possible pollution of other memories (instruction cache, branch prediction tables,...) can increase the execution time. Ignoring speculative execution might again lead to an erroneous classification of instructions (branches or memory accesses) in the global low-level analysis step, which may result in an underestimated WCET.

# 4. Towards a safe WCET estimation for processors with speculative execution

We have shown why the execution of the wrong path has to be carefully taken into account in order to obtain a safe WCET. In this section, we discuss the difficulties of modelling speculative execution as part of static WCET analysis.

The possible effects of speculative execution on the dynamic scheduling of instructions can probably be taken

into account without excessive complexity. For example, the delaying of the execution of later instructions due to the occupation of hardware resources by wrong path instructions could be included in the WCET estimation by systematically adding to the misprediction recovery penalty the longest functional unit latency. An other solution would consist in assuming in-order instead of out-of-order execution, but it would lead to a very pessimistic WCET estimation.

The effects of speculative execution on the content of memories may be harder to take into account within a purely static WCET analysis. We have seen that it can invalidate the results of the global low-level analysis: memory accesses (either to instructions or data) classified as "cache hits" might actually miss due to the pollution of the cache by the execution of the wrong path; branches classified as "well predicted" might actually be mispredicted due to the pollution of branch predictor tables. This means that, in order to produce a safe classification of instructions, the global low-level analysis should take into account the instructions of the wrong path. Now, we feel that considering wrong paths within static analysis is not straightforward, since it probably requires new algorithms for syntax tree or control-flow graph traversal. Moreover, the number of instructions or basic blocks to include in a wrong path depends on the processor state (occupancy of hardware resources) and can only be determined during the local low-level analysis. Thus, it appears that a correct modelling of speculative execution would require a very close interaction between the high-level (flow) and lowlevel (timing) analyses, which are usually carried out independently.

While decoupling the analyses of different components (caches, branch predictor, pipeline, ...) probably makes static WCET computation feasible, we wonder if the requirement of more interaction between these analyses to be able to take into account more and more advanced hardware features can be satisfied. If not, growing emphasis should be put on dynamic measurement (on real systems or simulators) to obtain accurate timing information while the static part of WCET estimation would focus on selecting the execution paths to explore with the goal of minimizing the measurement requirements.

#### 5. Conclusion

A lot of work has been done these last ten years to allow static estimation of the WCET for processors with advanced features like cache memories, pipelined execution, branch prediction. The most recent dynamically-scheduled processors implement speculative execution: when a branch instruction is predicted, the instructions belonging to the predicted path can be executed before that the branch is resolved. In this paper, we discussed the possible effects of executing the wrong path whenever a branch is mispredicted.

We have shown that wrong path execution can modify the scheduling of the correct path instructions and/or change the content of memories (instruction and data caches, branch predictor tables, ...). Then we have explained why ignoring these effects in WCET analysis can lead to an underestimated WCET, which can be dramatic for hard real-time systems.

We feel that usual static WCET computation techniques cannot accurately take speculative execution into account, since it would require a too complex interaction between global and local low-level analysis.

The next step of our work will consist in quantifying the effects described in this paper, with the goal of evaluating the probability that speculative execution increases the WCET. This measure will indicate whether new methods are to be investigated to take speculative execution into account within WCET analysis.

#### 6. References

- [1] A. Colin, I. Puaut, "Worst-Case Execution Time Analysis for a Processor with Branch Prediction", *Real-Time Systems*, 18(2):249-274, May 2000.
- [2] S. Jourdan, T.-H. Hsing, J. Stark, Y. Patt, "The Effects of Mispredicted-Path Execution on Branch Prediction Structures", Int. Conf. On Parallel Architectures and Compilation Techniques, Octobre 1996.
- [3] T. Lundqvist, P. Stenström, "Timing Anomalies in Dynamically Scheduled Processors", 20<sup>th</sup> IEEE Real-Time Systems Symposium, December 1999.
- [4] J. Pierce, T. Mudge, "Wrong-Path Instruction Prefetching", IEEE Int. Symp. On Microarchitecture, December 1996.

## Session V: Issues in WCET Analysis

#### Niklas Holsti

This last session of the workshop provided three presentations followed by discussion.

### Presentations

In the first presentation, Jörn Schneider noted that the execution time of an RTOS call often depends greatly on the application-defined context such as parameter values. Moreover, the execution time of application code can depend on the context set up by earlier RTOS calls, for example to disable preemption. These interactions decrease the accuracy of a WCET analysis done separately for the application and for the RTOS. Schneider proposed a combined WCET-and-schedulability analysis that considers such "meta-state information". Similar problems arise for software libraries and their clients ("library syndrome").

In the second presentation, Stefan Petters explained how statistical analysis gives a WCET estimate that is not an absolute upper bound but for which we know the risk of exceeding the estimate. If this risk is small enough in relation to other causes of system failure, the whole system is shown to be reliable enough. The approach resembles HW failure analysis and uses extreme value statistics.

Peter Puschner's final presentation suggested a radical change in HW and SW design to simplify timing analysis. Speculative HW mechanisms such as caches are replaced by deterministic ones such as SW-controlled prefetching. Code transformations reduce the number of SW execution paths to one. The time for this single path can be found by measurement or simple static analysis.

#### Discussion

The presenters listed several questions for discussion. Schneider asked for more problems in application-RTOS interaction; how WCET researchers approach the library syndrome and serve cache-sensitive, pipeline-sensitive and RTOS-aware schedulability analysis; and why industry is still measuring execution times. Petters asked if less than 100% WCET guarantee is acceptable; whether execution times can reasonably be described statistically; and what other sources of proof could enlarge confidence in a statistically derived failure rate. Puschner wanted to discuss the speed/quality trade-off in WCET analysis; whether current analysis is good and fast enough; and whether alternative approaches are needed, as his presentation suggested. Since this was the last session, the discussion also touched on general topics. The summary below is organized by issue rather than time order.

Both measurement and analysis are needed. Analysis is hampered by the need to build a model, while measurement uses the real system, which (as Bernat remarked) is its own best model. The misshapen mirror of the Hubble Space Telescope was brought up as an example of a costly failure that could have been prevented by measurement (Hillary). However,

the problem here was a faulty measurement (Vardanega). This risk also exists in measuring execution times, since SW engineers do not understand that measurement is like testing: it cannot show the *absence* of errors (Schneider) and may miss the critical case (Puschner). SW engineers should think about what they measure (Petters) and should combine measurement with analysis (Renaux) to validate the analysis (Schneider).

Vardanega remarked that typical applications use only a few RTOS services. In one example, only 27 of 256 services were used. This suggests that the interface between application and RTOS could be simplified, perhaps reducing the timing interactions as well.

The statistical approach presented by Petters was not much discussed. Hillary noted that it should fit in well with the reliability analysis already required for aircraft.

Puschner's proposals received a mixed reception. Petters said that the WCET researchers have no chance to influence processor design. Puschner acknowledged that the current designs will be continued for a decade, but called for optimism and suggestions for solutions. Holsti suggested that Puschner's approach could work with options added to current processors, while Bernat felt that the present pessimism of WCET results might simply be accepted. Vardanega stated strongly that the community should insist on open processor designs where the software engineer knows what will happen and how to model it.

Replying to Puschner's question, Gustafsson said that current WCET analysis is of course not good enough, but that we should not aim for "one best WCET analyser". Comparing WCET analysers to compilers, where there is a choice of hundreds of compilers, he felt that there will be a need for many different WCET analysers and methods, including the "single path" approach. However, Vardanega warned that compiler vendors are not seeing enough demand for WCET analysis to develop such tools within their R&D budgets.

Overall, the session showed that the research community is not ignoring practical problems, such as the application/RTOS interface, and is not stumped by the growing complexity of processors, but has at least two alternative approaches in addition to the straightforward (and still workable) approach of evolving more complex models.

## Why You Can't Analyze RTOSs without Considering Applications and Vice Versa\*

Jörn Schneider

Dept. of Computer Science, Saarland University, Germany E-mail: js@cs.uni-sb.de

#### Abstract

Traditionally worst case execution time (WCET) analysis tools are designed for the analysis of application code. The execution time of Real-Time Operating System (RTOS) services and the interaction between RTOS and application are usually not considered. When performing an RTOS aware schedulability analysis the WCETs of RTOS services are needed. At first sight the application of existing WCET analyzers on RTOS code should be straightforward and should deliver the same accuracy as for application code. The paper explains why this is not the case, and why the presence of an RTOS diminishes the accuracy of application code WCET analysis.

In addition to explaining why RTOSs should not be analyzed without considering application code and vice versa, the underlying problems are identified as well as enlightened by some examples and possible solutions are sketched. Eventually a comprehensive approach for WCET and schedulability analysis is proposed. The comprehensive approach solves almost all of the identified problems. Additionally, the synergetic approach opens the road towards very promising, novel possibilities, like automatic trouble spot identification.

#### 1 Introduction

To assure that the timing constraints of a hard real-time system are met it is of paramount importance to consider the temporal impact of all parts of the system. Provided that the temporal behavior of the environment is given, the main parts of many contemporary real-time systems are the hardware (e. g. a microprocessor with caches and pipelines), the real-time operating system (RTOS), and the application software. The research community addressed the problem from two sides in the past. Considering the real-time operating system behavior has been gradually introduced in continously improved schedulability analysis approaches. The temporal parameters of application software are predicted by worst case execution time (WCET) analysis. The hardware influence was initially captured by the WCET analysis part alone. For simple hardware there is no need to consider microarchitectural effects at the schedulability analysis level. More recent work considers the impact of modern microarchitectural features like caches and pipelines at the schedulability analysis level [13]. However, this did not close the last gap yet. One remaining gap is the computation of worst case execution times of RTOS services. Ideally, these times should be obtainable with the WCET analyzers designed for analyzing application code without further measures. Certainly this is possible, but at what costs? A recent study of Colin and Puaut exemplifies that the results of such an approach are poor [4].

This paper shows that analyzing the WCET of RTOSs without considering the application, inevitably leads to poor results. Additionally, it is shown that the presence of a multitasking RTOS diminishes the accuracy of application code WCET analysis. The underlying problems of these phenomena are identified and explained by examples. Furthermore, possible solutions are sketched. The found solutions point towards a direction that leads to the following observations: 1. Combining schedulability analysis and WCET analysis in a straightforward way means that these two techniques work independently but not jointly. 2. Deploying WCET analysis and schedulability analysis in a comprehensive framework leads to synergies that are a remedy for the diagnosed problems and open further very promising possibilities. Therefore, a comprehensive framework for WCET and schedulability analysis is proposed.

The paper is organized as follows. The next section presents related work from the fields of WCET and schedulability analysis. Section 3 identifies the problems encountered when analyzing the WCET of RTOSs and sketches possible solutions. The case of analyzing the WCET of applications in presence of a multitasking RTOS is discussed

<sup>\*</sup>This work has been partly supported by DFG (German Research Foundation), Transferbereich 14, and by the European Commission, ARTIST Project (Advanced Real-Time Systems Information Society Technologies).

in Section 4. Section 5 briefly describes the proposed comprehensive analysis framework. Section 6 provides conclusions. The last section reflects some related results of the discussions at the WCET workshop.

#### 2 Related work

Current WCET analyzers aim at analyzing application code. The execution time of RTOS services and the interaction between RTOS and application are usually not considered. However, they are capable to consider the timing impact of modern microarchitectural features. The work of Ferdinand [5] shows how to consider the impact of setassociative instruction and data caches. Later joined work with the current author extends this approach to superscalar pipelines [14]. More recent work of the USES (University of the Saarland Embedded Systems) group augments the method to unified instruction and data caches [6] and dynamic branch prediction, speculative execution as well as out-of-order execution [16].

Colin and Puaut published a first study on program analysis based WCET estimation of RTOS primitives [4]. They analyze twelve system calls (including the scheduler) of the RTEMS kernel [11] for the Pentium I CPU and report several problems in applying their WCET analysis. The main problems were due to applying WCET analysis without considering schedulability analysis and application code. For instance the loop bound of the RTEMS scheduler could not be derived because it depends on the number of task arrivals during its execution (the scheduler loops until no further arrivals are noticed). The average WCET overestimation reported is 86%.

Besides the publications on WCET analysis the area of schedulability analysis, especially on fixed priority scheduling is related. Several publications in this area demonstrate how to consider many RTOS features within schedulability analysis. Katcher et al. provide utilization based schedulability tests for four different scheduler implementations [8], two for event driven and two for time driven scheduling (tick scheduling). The paper is restricted to non-communicating tasks in periodic task sets. Burns and Wellings show for the attitude and orbital control system of the Olympus satellite how features like context-switch times, sporadic tasks, system clock (tick scheduling of periodic tasks), and release jitter can be considered in response time analysis [2]. These publications tend to neglect the effect of microarchitectural features. Katcher et al. do not mention cache and pipeline behavior at all. Burns and Wellings switch off the cache and ignore pipeline effects by making pessimistic assumptions.

Other publications show how to consider microarchitectural preemption costs in schedulability analysis. Most of them incorporate cache-related preemption costs by adding

fixed cache reload costs for each preemption (with the exception of [13], which is the basis of the author's current work). Thereby, it is impossible to consider the compensation of cache misses by pipeline effects. Therefore, these methods are referred to as isolated methods. Pipeline effects are often completely ignored. Basumallick and Nilsen describe an extension of the utilization based rate monotonic analysis (RMA) of Liu and Layland [10], which considers cache-related preemption costs [1]. The WCET of a task is increased by the worst case cache refill costs that this task may impose on any preempted task. These additional costs can be either the size of the cache or the size of the memory area occupied by the preempted task. Busquets-Mataix et al. show that the RMA-method is worse than a comparable response time analysis method [3] because of its pessimistic assumptions on the schedulable workload. Their CRTA (cached version of RTA) [3] adds cache reload costs for each preemption of a task. The authors mention five ways to estimate this cache related interference. CRTA supports two of these five ways: either the costs for refilling the entire cache or the time to refill all cache lines possibly displaced by the preempting task can be used. Lee et al. [9] present a sophisticated isolated approach. The approach considers the phasing of tasks. For each program point of each task an upper bound on the worst case cacherelated preemption costs at this program point is calculated. These results are used in a second phase to derive an upper bound on the overall worst case cache related preemption costs during the response time of the preempted task. The second phase uses a linear programming technique that executes in each iteration of the response time analysis. As in [13] the relationship between the preempting task and the set of tasks that is possibly running when the preemption occurs is taken into account.

All three approaches [1, 3, 9] are isolated methods and do not consider pipeline-related preemption costs. In contrast to them the isolated method presented in [13] considers pipeline-related preemption costs. No isolated approach is able to take into account the overlapping of preemption caused cache effects with pipeline hazards. To overcome this drawback in [13] an integrated method was presented. The integrated method considers cache-related preemption costs implicitly and pipeline-related preemption costs explicitly. Among these approaches only [9] and [13] partly considered the influence of RTOSs by regarding the scheduler overhead.

#### **3** Analyzing RTOSs

#### **3.1** What are the problems?

The WCET of RTOS services is highly dependent on the application using them. Table 1 gives a systematic list of

such dependences. Examples of such dependences can for

Determining factors of the WCET of RTOS services	Examples
Non-constant call parameters in application code	Any service with call parameter dependent control flow
Static, application dependent configuration parameters	Any service with loop bounds depending on no. of RTOS objects (e. g. tasks, resources)
Cache state	Replacement of RTOS owned cache sets by application code
Calling history of RTOS services	Scheduler execution after disabling preemptions
Calling context	Call to system service from task, interrupt or operating system level

Table 1. Sources of WCET variations of RTOS services.

instance be found in the RTEMS code, and in the code of osCAN (an OSEK [12] implementation by Vector Informatik).

#### **3.2** How can these problems be addressed?

**Non-constant calling parameters** When system calls are analyzed as part of the application, any knowledge about parameter values (e. g. obtained by a value analysis [15, 6]) can be used to derive sharper bounds on the WCET.

**Static configuration parameters** The configuration parameters (e.g. number of tasks and memory mapping of tasks) are fixed for a particular application. Therefore, they can be considered either manually or automatically by WCET analysis as well as schedulability analysis.

**Cache state** If the cache is not partitioned in a special way, application code or data might displace cache sets occupied by the RTOS. Therefore, no isolated WCET analysis of the RTOS can benefit from positive cache effects caused by previous runs of RTOS services. It might even be impossible to consider the positive intrinsic cache effects of RTOS services. RTOSs are usually designed to minimize

the number and duration of non-interruptible code sections. It is impossible for an isolated analysis to predict the negative impact of application interrupts outside these few code sections, unless all positive cache effects are ignored. In the case of a combined analysis it is possible to bound the effect of application caused cache replacements as it has been shown in [13] for the application analysis. This is reached by considering the preemption related cache effects within the WCET analysis. Whether this approach delivers better results than traditional ones depends on the application. For a high number of preemptions the novel approach is superior, for systems with a low workload the traditional methods are better (cf. [13]). In [13] it is also shown that for certain modern microprocessor types this approach is imperative in order to avoid underestimations of response times.

**Calling history of RTOS services** It not only affects the WCET of RTOS services, but often has an immediate impact on the task response time as well (e.g. an RTOS service called to disable preemption eliminates the subsequent interference by other tasks). A good schedulability analysis should consider these effects. Therefore, the history information should be statically predicted anyway and can also be used by WCET analysis.

**Calling context** When using a combined analysis, the calling context can easily be regarded. The WCET analysis can for instance ignore paths that become infeasible because of the specific calling context.

#### 4 Analyzing applications

This section discusses problems arising in presence of multitasking RTOSs. The two subsections treat the problems in the same order. First the problem domain of data values is considered. Issues that arise due to isolated analysis of application and RTOS code come second. Not all of these issues can be addressed by a mere integration of application and RTOS WCET analysis. The last parts of either subsection and Section 5 cope with this enigma.

#### 4.1 What are the problems?

The data values used in application code can play a large role in computing the WCET.<sup>1</sup> Examples are: loop bounds, addresses of memory references and infeasible paths. For the WCET analysis to profit from this fact a static prediction of value ranges is necessary. The value analysis described in [15, 6] provides this functionality for instance.

<sup>&</sup>lt;sup>1</sup>This holds for RTOS code also. Nevertheless the subject is discussed in this section because that is where WCET analysis comes from and because applications are usually more data-driven than RTOSs.

However, tools of this kind are—like any available WCET tools—designed for sequential programs. The following issues arise in presence of a multitasking RTOS:

- **Shared application memory** Accesses by other tasks may change the value of data in such areas.
- **RTOS data structures** Any RTOS data structure not unique to the analyzed task might be changed by RTOS services called in other tasks. Even data structures unique to a task might be manipulated by other (user or RTOS) programs.
- **Memory mapped I/O** The values read from those areas are mainly determined by the environment and accesses are non-cachable.

The WCET analysis of application code should consider the WCET of the system calls used. A seemingly attractive approach is to initially ignore the system calls within the application WCET analysis and thereafter add system call WCETs obtained by an isolated analysis of the RTOS. However, there are good reasons not to do so (see Table 2).

No.	Problem description
1	RTOS WCETs are systematically overestimated
	(shown in Section 3)
2	Information about correlation of worst case paths
	and number+context of RTOS calls is destroyed
	$\Rightarrow$ only a pessimistic approach can still deliver
3	Cache and pipeline effects caused by RTOS calls
	cannot be considered in application WCET
4	It is impossible to consider positive effects of
	concepts existing only in presence of multi-
	tasking RTOSs, for instance RTOS calls
	dynamically raising the application priority
	(e. g. by disabiling preemption or interrupts, or by occupying resources)
	or by becapying resources)

## Table 2. Problems of analyzing applicationsisolated from the RTOS.

There is a significant difference between problem descriptions 1 through 3 and the classes of problems alluded to by description 4 of Table 2. The former difficulties occur also together with the isolated WCET analysis of library functions, the latter not.

#### 4.2 How can these problems be addressed?

Because of shared application memory and RTOS data structures, a value analysis has either to ignore such data, or has to be enhanced to a multi-task-analysis. The latter is not trivial. Memory mapped I/O areas have to be excluded from the value analysis since they are volatile.

Section 3 shows that the WCET of RTOS services depends on the call situation. This call situation subsumes the factors given in Table 1. Some aspects of the call situation are not unique to applications running on RTOSs. These aspects can be identified already, when stand-alone applications with calls to library routines are considered (one could replace the word RTOS with library in the problem descriptions number 1 through 3 of Table 2 and the statements would still be true to some extent). Aspects like these can be addressed by embedding the analysis of library/system calls within the application WCET analysis. The embedding can be implicitly or explicitly. Embedding implicitly means that the calls are treated like ordinary function calls. The input data structure (e.g. the control flow graph) of the WCET analyzer has to contain all needed information to analyze such calls. If the WCET analyzer uses machine code as input (e.g. the one described in [6]), this can even be done without providing the user with the library/RTOS source code. Embedding explicitly means that two independent WCET analyzers (or two instances of the same analyzer) are used, one for the application and one for the RTOS. The RTOS WCET analyzer can be a black box that takes the code of the RTOS service as well as collected information about calling parameters, static configuration parameters, cache state and calling context as input (see Subsection 3.2).

However, introducing an RTOS in the considered scenario adds completely new qualities to the problem of WCET analysis (represented by item no. 4 of Table 2). These are issues that cannot be addressed by a mere integration of application and RTOS WCET analysis. Rather a high-level view is needed to consider them in WCET and schedulability analysis. Several such high-level concepts can be identified that co-determine the temporal behavior of the tasks of an RTOS-based system. These concepts are for instance the effective priority of tasks, the RTOS mode (e.g. initialization or normal operation mode), application modes, task states, and the system level (task, interrupt or RTOS level). Those high-level concepts have a certain meaning when viewing the system as a whole rather than as a bunch of independent programs at a microarchitectural level. At run-time the properties of these concepts have a defined state at each point in time. We define the meta-state of a task to be the set of these states. In a static approach, at best partial knowledge of the meta-state of a task can be obtained. To address the RTOS specific problem domain, partial knowledge can be collected that allows

to compute the worst case response time of tasks more accurately. This includes exploiting meta-state information to compute sharper bounds on WCETs of tasks as well as sharper bounds on microarchitecture-related preemption costs.

## 5 Proposal for a comprehensive WCET and schedulability analysis approach

The authors present work on a comprehensive WCET and schedulability analysis approach exploits meta-state information to compute sharper bounds on WCETs of jobs (tasks and interrupt service routines) and interesting code sections and on microarchitecture-related preemption costs. The framework exploits the following aspects of the meta-state of a job: effective priority of a job (determined by: locked interrupts, preemption lock, occupied resources), RTOS mode and current system level. The meta-state information is exploited as follows:

- 1. Extrinsic cache effects are considered by the cache analysis (which is a part of the WCET analysis) in dependence of the effective priority at each program point of the analyzed job.
- Pipeline-related preemption costs are individually computed for each job, again in dependence of the effective priority, and are considered during the schedulability analysis.
- 3. The WCET of jobs and code sections is computed for the proper RTOS mode (initialization mode or normal operation mode) and for each RTOS mode a separate schedulability analysis is undertaken.
- 4. The current system level is considered when the WCET of system calls is computed.

Similar to the cache- and pipeline-sensitive schedulability analysis described in [13] the cache-related preemption costs are incorporated in the WCET while the pipelinerelated preemption costs are explicitly considered during the schedulability analysis.

The above sketched framework uses the WCET analysis tool described in [6]. The WCET analyzer is loosely coupled with the surrounding tools. It is guided with the help of the obtained meta-state information in order to compute sharper bounds on the WCET and the microarchitecturerelated preemption costs. A detailed explanation of this method is beyond the scope of this paper.

The described approach yields a valuable side effect. With only little more effort information highly esteemed by system developers can be derived. For instance, the individual latency of each interrupt together with the causing code section(s), code sequences dominating the WCET of jobs due to destroyed cache contents and the causing memory references of higher priority tasks, and critical code sections dominating blocking time of tasks. This opens the door towards precisely aimed optimization efforts to make infeasible task sets schedulable or provide "space" for additional tasks, without requiring expert knowledge.

#### 6 Conclusion

The paper showed that analyzing WCETs of RTOS-based real-time systems—whether of RTOS services or of application code—requires other than the established approaches. The underlying problems were explained by examples and classified. Additionally, it was sketched how the individual problems can be addressed. Finally a comprehensive approach for WCET and schedulability analysis was proposed. The proposed approach shows how it is possible to overcome most of the obstacles obstructing the path toward comparative results of WCET analysis for RTOSbased systems.

As far as the author knows this is the first proposal to consider static computable information about the metastate of jobs. This information can be used to guide the WCET analysis of jobs and interesting code sections to achieve more accurate WCETs and to sharpen the results of the analysis of microarchitecture-related preemption costs. Therefore, the proposed approach has a much higher potential for good results than any combination of traditional methods, which consider either RTOS or microarchitecture related costs.

The comprehensive method yields a valuable side effect. With only little more effort information highly esteemed by system developers can be derived. For instance the individual latency of each interrupt together with the causing code section(s), code sequences dominating the WCET of jobs due to destroyed cache contents and the causing memory references of higher priority tasks, and critical code sections dominating blocking time of tasks. This opens the door towards precisely aimed optimization efforts to make infeasible task sets schedulable or provide "space" for additional tasks, without requiring expert knowledge.

#### 7 Discussion Results

Several related topics were discussed at the WCET Workshop in Vienna. Some of them raised by the workshop version of this paper. The industrial representatives appreciated especially the comprehensive approach and acknowledged its benefit for industrial practice. The growing importance of schedulability analysis solutions has been emphasized by the industrial as well as the research community.

The infeasibility of the traditional approaches, especially measuring in presence of caches and pipelines [7], has been acknowledged by the industrial representatives. The steady progress of microprocessors with caches and complicated pipelines into the real-time area has been reported from many participants. Together this means that at the moment no solution exists that allows to guarantee hard real-time behavior and to benefit from modern hardware in the presence of RTOSs. Therefore, the proposed comprehensive framework is urgently needed. The surplus of the solution, e. g. the possibility to identify trouble spots and thus allow precisely aimed optimizations, was appreciated especially by the industrial representatives.

#### Acknowledgements

Many members of the compiler design group at the University of the Saarland, especially the members of the USES (University of the Saarland Embedded Systems) group, deserve acknowledgement. Reinhard Wilhelm, Daniel Kästner, and Stephan Diehl carefully read draft versions of this work and provided many valuable hints and suggestions.

I want to express my gratitude to the company Vector Informatik for making available their implementation of the OSEK real-time operating system standard.

I would like to thank the anonymous reviewers for their helpful comments and the workshop participants for their comments during insightful discussions.

#### References

- S. Basumallick and K. Nilsen. Cache Issues in Real-Time Systems. ACM SIGPLAN Workshop on Language, Compiler and Tool Support for Real-Time Systems, 1994.
- [2] A. Burns and A. Wellings. Engineering a Hard Real-time System: From Theory to Practice. *Software—Practice and Experience*, 25(7):705–726, 1995.
- [3] J. V. Busquets-Mataix, J. J. Serrano, R. Ors, P. Gil, and A. Wellings. Adding Instruction Cache Effect to Schedulability Analysis of Preemptive Real-Time Systems. In *Proceedings of the IEEE Real-Time Technology and Applications Symposium*, pages 204–212, June 1996.
- [4] A. Colin and I. Puaut. Worst-Case Execution Time Analysis of the RTEMS Real-Time Operating System. In Proceedings of the 13th Euromicro Conference on Real-Time Systems, pages 191–198, Delft, The Netherlands, June 2001.
- [5] C. Ferdinand. Cache Behavior Prediction for Real-Time Systems. Dissertation, Universität des Saarlandes, Sept. 1997.
- [6] C. Ferdinand, R. Heckmann, M. Langenbach, F. Martin, M. Schmidt, H. Theiling, S. Thesing, and R. Wilhelm. Reliable and Precise WCET Determination for a Real-Life Processor. In *Embedded Software Workshop*, Lake Tahoe, USA, Oct. 2001.

- [7] N. Hillary and K. Madsen. You Can't Control what you Can't Measure, or Why it's Close to Impossible to Guarantee Real-Time Software Performance on a CPU with On-Chip Cache. In Proceedings of the WCET Workshop of the 14th Euromicro Conference on Real-Time Systems, Vienna, Austria, June 2002.
- [8] D. I. Katcher, H. Arakawa, and J. K. Strosnider. Engineering and Analysis of Fixed Priority Schedulers. *IEEE Transactions on Software Engineering*, 19(9):920–934, 1993.
- [9] C.-G. Lee, J. Hahn, Y.-M. Seo, S. L. Min, R. Ha, S. Hong, C. Y. Park, M. Lee, and C. S. Kim. Enhanced Analysis of Cache-Related Preemption Delay in Fixed-Priority Preemptive Scheduling. In *Proceedings of the IEEE Real-Time Systems Symposium*, pages 187–198, Dec. 1997.
- [10] C. L. Liu and J. W. Layland. Scheduling Algorithms for Multiprogramming in a Hard Real-Time Environment. *Journal of the ACM*, 20(1):46–61, 1973.
- [11] On-Line Applications Research Corporation, Huntsville, AL, USA. RTEMS Applications C User's Guide. Edition 4.0, Oct 1998. http://www.oarcorp.com/RTEMS/ rtems.html.
- [12] OSEK/VDX Open systems and the corresponding interfaces for automotive electronics. OSEK/VDX Operating System. Version 2.2, Sept. 2001. http://www. osek-vdx.org.
- [13] J. Schneider. Cache and Pipeline Sensitive Fixed Priority Scheduling for Preemptive Real-Time Systems. In Proceedings of the 21st IEEE Real-Time Systems Symposium 2000, pages 195–204, Nov. 2000.
- [14] J. Schneider and C. Ferdinand. Pipeline Behavior Prediction for Superscalar Processors by Abstract Interpretation. In Proceedings of the ACM SIGPLAN Workshop on Languages, Compilers and Tools for Embedded Systems, volume 34 of ACM SIGPLAN Notices, pages 35–44, May 1999.
- [15] M. Sicks. Adreßbestimmung zur Vorhersage des Verhaltens von Daten-Caches. Diplomarbeit, Universität des Saarlandes, Nov. 1997.
- [16] S. Thesing, M. Langenbach, and R. Heckmann. Pipeline Behavior Analysis for Real-Time Systems by Pipeline Modeling. In Proceedings of the Work-in-Progress Session of the 14th Euromicro Conference on Real-Time Systems, pages 73–76, Vienna, Austria, June 2002.

## How much Worst Case is Needed in WCET Estimation? \*

Stefan M. Petters Department of Computer Science University of York United Kingdom Stefan.Petters@cs.york.ac.uk

### Abstract

Probabilistic methods provide probability density or mass distributions for the execution time or the assumed worst case execution time instead of a single WCET value. While the resulting probability tends to fall towards zero quickly, the actual zero value, i.e. the 100 % guarantee, is reached only with unreasonable overestimation of the real WCET. In order to cope with this, this paper proposes to use similar techniques to hardware dependability analysis, where a 100 % guarantee is physically impossible and a certain, usually very small amount of risk is acceptable.

## 1 Motivation

Modern high performance processors include many features which usually make cycle true simulation infeasible, impose impractical limitations on the code or operating system, or the methods used to capture these effects have to introduce simplifications that lead to results which may be up to an order of magnitude beyond the physical possible WCET.

One possibility to get around this problem is the deployment of statistical methods. Throughout this paper the validity of these methods is assumed. Additionally it has to be assumed that the methods provide a description of the program behaviour which correctly bounds the execution time for all modi of operation. This has to be especially considered for systems in *red button mode*, where this critical mode has to receive seperate and close consideration. A major problem of such approaches is that they result in approximations of the (worst case) execution time, whose probabilities are non-zero, but very small for a long way beyond the physical WCET. The question now is, whether one has to go for the zero-probability of an error, which tends to be as pessimistic as the simplified WCET analysis, or if a probabilistic guarantee suffices. While the first case, gains no real advantage, the second has an open issue, which probabilistic guarantee to accept as good.

## 2 Probabilistic WCET Analysis

Research in probabilistic WCET analysis can be divided in two categories:

- Approaches using observed test cases to reason about the probability of an execution time not observed during the tests.
- Approaches analysing small parts of the program in order to reason about the probability of different combinations of the results of the smaller units.

<sup>\*</sup> The work presented in this paper is supported by the *European Union* as part of a research programme on "Next TTA"

As there are currently no publications in the second area, we will focus on an example of the first reserach area. Stewart Edgar uses a black box approach in [1]. The description of the method here can only be very coarse and the reader should have a look at the original papers on this topic (e.g. [1, 2]).

The program is run several times, with random input data and the end-to-end execution time of the program runs are measured. As it is obvious, the measurements will not likely cover the physical WCET of the program on that processor, extreme value statistics are deployed to reason about the execution time longer than any experienced during measurement. Extreme value statistics are concerned with modelling the left and/or right hand tail of a probability distribution, as opposed to the modelling of the average case with conventional statistics. This induces that outliers in the measurement data, which are usually disregarded with conventional statistics, have considerable impact on the modelling parameters of extreme value probability density functions.

Extreme value statistics are well known in the area of financial risk assessment and civil engineering. In the latter case the assessment of maximum wind speeds or flood levels is computed utilising these technique in order to dimension the statics of buildings. There are three type of extreme value probability density functions, described with a theorem which corresponds to the central limit theorem of the normal distributions. As a necessary precondition to apply this technique, the underlying random variables have to be independent and identically distributed.

For the approach the most simple solution of a *Gumble distribution* has been chosen. This model only uses deviation and mean of the random variable, in our case the observed execution time. The following equation show the Gumble probability density function and the cumulative Gumble probability density function:

$$g(t) = exp\left(-\left(e^{-\frac{t-\mu}{\sigma}} + \frac{t-\mu}{\sigma}\right)\right)\frac{\mu}{\sigma} \quad (1)$$
  

$$G(t) = exp\left(-e^{-\frac{t-\mu}{\sigma}}\right) \quad (2)$$



Figure 1: Sample Measurement Data and Extreme Value Approximation.

The cumulative variant expresses the probability of an execution time below the value t.

An example execution time measured and the corresponding Gumble distribution is given in figure 1. The measured times are given in a kernel density transformed representation. The transformation is used to display discrete data as a continuous curve and thus allowing the comparison by inspection with the extreme value approximation.

A major drawback of using the Gumble distribution to approximation is the non-zero probability for execution times, except for  $\pm\infty$ . While the probability of the execution time exceeding  $20\sigma$  beyond the mean is only  $2.06E^{-9}$ , the "risk" is still there. As experiments show, this probability reaches quickly  $1.0E^{-20}$  and less with an overestimation of some 10% (cf. [3]). As an important aspect it has to be noted, that the probabilities targeted in this paper, correspond to *execution times* beyond anything observed during the tests.



Figure 2: Typical Variation of Fault Ratio of Hardware Components over Time [4].

### **3** Hardware Considerations

This section will give a short introduction in the mechanisms to risk assessment of hardware components. Figure 2 shows the typical distribution of hardware faults in electronic equipment over time.

During the period of *burn in* the probability of hardware failure is higher, due to faults in the productions of the components. A good example for such a behaviour are errors due to the statistic deviation in the doting of semiconductors. To avoid the high probability of failures in the burn in period, the components are in general case run for a time before deployment in a dependable system to weed out bad components. This process is in most cases sped up by undertaking this testing phase under more extreme circumstances than the system has to endure in real operation (e.g. heat, cold, mechanical stress). Thus a production error that might show up only after months or years down the line is uncovered after a few hours or days of operation.

After the burn in time, the hardware components reach a more or less constant failure rate of  $\lambda$ . Usually this useful lifetime is quite long. In the end the wear out sets in, where, for example, saturation effects<sup>1</sup> in the semiconductor set in. The failure rate  $\lambda$  after burn in as well as the average life time of a given hardware component is usually known. The reliability R(t) of a component not to fail is given in equation 3.

$$R(t) = e^{-\lambda t} \tag{3}$$

For the computation of system failure usually the *mean* time to failure (MTTF) (i.e.  $\lambda$ ) is taken to compute the overall systems failure rate. Since the usual failure rate is less the one failure in the lifetime  $T_{\text{life}}$  of a product, the failure rate can be transformed into a failure probability for the lifetime  $p_{\text{life}}$  of the system. This failure probability can be computed using equation 4.

$$p_{\text{life}} = \int_{0}^{T_{\text{life}}} R(t) dt \qquad (4)$$

A similar reasoning may also be applied to software components. The major difference between software and hardware components is the discrete nature of failures of the software components as opposed to the continuous nature of failures of the hardware components. Assuming the program has no algorithmic errors, exceeding a computation time alloted to the program can be considered a software failure in real-time systems. A basic necessity for this is the assumption that the probability for an overrun of the alloted time for an individual run  $p_{\text{excess}}$  is known and constant for all runs. Additionally the maximum amount of task releases for any given time is essential for the computation. This is usually defined as a minimal inter arrival time  $T_{\text{task}}$ 

The probability of a failure over the lifetime of the product is computed using equation 5.

$$p_{\text{life}} = 1 - (1 - p_{\text{excess}})^{\frac{T_{\text{life}}}{T_{\text{task}}}}$$
 (5)

Defining an acceptable failure probability during the lifetime, which would be in the order of magnitude of

<sup>&</sup>lt;sup>1</sup>One problem in the semiconductor industry is that the doting of

semiconductors may be done by diffusion and these donated atoms tend to start drifting inside the semiconductor.

the failure probability of an hardware failure, it is easy [4] N. Storey, Safety-Critical Computer Systems. to compute an acceptable  $p_{excess}$  transforming equation 5 into:

Addison–Wesley Publishing Company, 1996.

$$p_{\text{excess}} = 1 - p_{\text{life}}^{\frac{T_{\text{task}}}{T_{\text{life}}}}$$
 (6)

#### 4 Conclusion

While the number of publications in the area of probabilistic WCET estimation is quite limited up to now, the number of people working on this issue is becoming larger. Interpreting an overrun of an assumed value for the WCET of a program as a software fault, similar probabilistic techniques as for hardware component failures may be used. This is particular useful whenever probabilistic methods are utilised to reason about the WCET, as these methods tend to provide probability density or mass distributions to describe the WCET instead of a single value. Special care has to be taken for emergency scenarios, as a failure in such a scenario is usually less acceptable than during normal execution. The validity and applicability of this method is subject to discussion.

#### References

- [1] A. Burns and S. Edgar, "Statistical analysis of WCET for scheduling," in Proc. of the IEEE Real-Time Systems Symposium (RTSS'01), (London, United Kingdom), Dec. 4-6 2001.
- [2] A. Burns and S. Edgar, "Predicting computation time for advanced processor architectures," in Proceedings of the 12th Euromicro Conference on Real-Time Systems, (Stockholm, Sweden), June 19-21 2000.
- [3] S. M. Petters, Worst Case Execution Time Estimation for Advanced Processor Architectures. PhD thesis, Institute of Real-Time Computer Systems, Technische Universität München, Munich, Germany, 2002.

## Is Worst-Case Execution-Time Analysis a Non-Problem? — Towards New Software and Hardware Architectures \*

Peter Puschner Institut für Technische Informatik Technische Universität Wien A1040 Wien, Austria E-mail: peter@vmars.tuwien.ac.at

#### Abstract

Despite the scientific advances in the research area of worst-case execution-time (WCET) analysis, there is hardly any industrial impact of the research solutions presented so far. This seems to be due to the high complexity of implementing and using the proposed WCET approaches.

This paper discusses what makes WCET analysis complex and proposes to use adequate hardware and software architectures to improve the predictability of program timing, thus simplifying WCET analysis.

#### 1 Introduction

Research in worst-case execution-time (WCET) analysis has been around for one and a half decades. During this period a number of different approaches to WCET analysis, including solutions for modelling hardware features and characterizing possible execution paths of real-time tasks have been found [5]. Still, the results of WCET-analysis research have hardly any impact on the industrial practice of timing analysis. This seems to be due to the high complexity of the implementation and use of WCET analysis tools. In addition, WCET research always seems to lag one step behind the advances in micro-processor technology – whenever WCET research manages to deal with the features of one hardware generation the next genera-

tion of processor and hardware architectures, equipped with novel speedup features, are already there.

This paper proposes to use new, adequate hardware and software architectures to improve the temporal predictability of programs, and thus reduce the complexity of WCET analysis. Hardware architectures for future real-time systems must allow to determine instruction execution times locally by inspecting single instructions and only a small number of instructions preceding them. Software architects have to investigate into programming techniques that reduce the number of input-data dependent branching decisions in the software, thus reducing the number of different execution paths of a program.

The further sections present our considerations in more detail and propose possible solutions. Section 2 argues about the main issues of WCET analysis and explains why the analysis is complex. Section 3 presents two possible ways to reduce the complexity of WCET analysis, one using special-purpose hardware features and the other based on the so-called singlepath programming paradigm. Section 4 concludes the paper.

### 2 The Complexity Dilemma of WCET Analysis

There is no doubt that WCET analysis as it is currently used is a complex problem. It has been shown that, in general, the number of paths to be analyzed for an exact WCET analysis of a piece of code grows exponentially with the number of consecutive branches in the control flow of the analyzed code. This state-

<sup>\*</sup>This work has been supported by the IST research project "High-Confidence Architecture for Distributed Control Applications (NEXT TTA)" under contract IST-2001-32111.

ment assumes that the code (a) is coded in traditional style (i.e., not applying programming techniques that focus on ease of WCET prediction) and (b) is to be executed on a modern high-performance processor architecture that includes caches and pipelines. Except for very simple programs this high complexity makes the full path enumeration needed for an exact WCET analysis intractable [3].

Current approaches to WCET analysis deal with this complexity in two ways:

- Calculate a high-quality WCET bound by accepting long computation times for the analysis.
- Trade the feasibility or speed of analysis for quality, by using simplified but pessimistic models of the possible software behaviours and the hardware timing.

The dilemma of WCET analysis is that neither of these approaches is acceptable in a commercial setting. On the other hand, using current hardware and software architectures does not allow for a better solution – the complexity of the problem simply is there. This raises the question if the current approaches to WCETanalysis are the correct answer to the problem of task timing analysis, or if current WCET research is focussing on the wrong problem.

#### 2.1 Sources of Complexity

Puschner and Burns [5] identified two central factors that determine the WCET of a program in a given application context:

- 1. the possible sequences of program actions in a given application, and
- 2. the time needed for each action in each of these possible sequences.

Clearly, both factors do not only determine the WCET of the code, but also the complexity of WCET analysis. Possible sequences of actions (instructions) depend on the algorithm that has been chosen to implement a solution to a problem and the code manipulations the compiler performs during compilation. The time needed for each action (instruction) depends on the features and configuration of the machine (hardware) on which the actions are executed. A number of principles applied in typical modern hardware and code design can be made responsible for WCET complexity. In the following we focus on two such principles, one for hardware and one for software.

Hardware Speedup by Speculation: This is the principle found in hierarchical memory architectures, e.g., cache. Instructions or data are loaded into (and kept in) small buffers (caches) with short access times these access times are typically much shorter than the access times of the larger store that holds larger portions of instructions and data — with the intention to speed up future memory accesses. The decisions about which items are to be loaded, kept, and replaced in cache are usually guided by heuristics, i.e., speculation about which items might be accessed in the near future.

The use of speculative decision mechanisms leads to variable memory access times. The duration of each particular memory access, in turn, depends on the state in which the preceding (and potentially very long) sequence of operations have left the cache. Both effects (the fact that memory access times vary and the dependency of actual memory access times on the execution history) taken together contribute to the complexity of WCET analysis.

Software Optimization for Frequent Scenarios: Realtime programmers use algorithms and programming techniques that have proven to be effective for non real-time applications. In non real-time applications, speed optimization for the most probable (i.e., frequent) scenarios is the primary goal. Temporal predictability is not an issue. In order to favour frequent cases, non real-time algorithms choose the actions to be performed based on input data. Input-data dependent control decisions, however, cause programs to execute on different execution paths with different execution times. As a consequence the number of different cases to be considered by the WCET analysis is potentially high.

#### **3** Possible Ways Out of the WCET Dilemma

This section illustrates the potential of alternative hardware and software architectures to simplify WCET analysis significantly. It provides alternatives to each of the two mentioned design principles.

#### 3.1 Hardware Speedup: Control Instead of Speculation

In contrast to non real-time applications, (hard) real-time applications primarily require temporal predictability. Appropriate hardware designs therefore support WCET analysis via predictability. This can be achieved by using memory hierarchies that exercise absolute control on the contents of fast buffers instead of relying on speculation. Rather than hoping that future memory accesses result in a cache hit, adequate prefetching strategies make the contents and thus access times of high-speed memory easy to predict. A memory architecture that achieves predictability by prefetching has been proposed a number of years ago [2]. Unfortunately, alternative memory architectures have not been further explored to this date.

#### 3.2 Software: Getting Rid of Input-data Dependencies

The second problem we mentioned is that traditional algorithm design and optimization yields code that behaves differently for different input data. To circumvent this problem and allow for a simple analysis, program behaviour must be less dependent on inputdata values. By reducing input-data dependencies the number of paths to be considered during WCET analysis gets smaller and, as a consequence, the complexity of the analysis decreases.

Following this concept we developed the singlepath paradigm [6]. The single-path paradigm yields programs that are fully temporally predictable. The central idea of the paradigm is to generate programs whose behaviour is completely independent of input data and which thus always execute on the one and only possible execution path.

Single-path programming builds upon a code transformation that removes data-dependent branching statements from the code. This code transformation is capable of transforming every WCET-analyzable piece of code into code with a single path. The transformation uses two different strategies to convert statements with if-then-else and loop semantics, respectively. Ifthen-else and other sequential branching statements with an input-data dependent branching condition are transformed into strictly sequential code by using *if*- *conversion*, [1]. Loops with input-data dependent termination are replaced by new loops with a constant the maximum — iteration count. After an input-data dpendent loop construct has been replaced, the conversion incorporates the original termination condition of the loop into the condition of a new *if* statement that it places around the body of the loop. As a last step, if-conversion is applied again to eliminate the newly generated *if* statement from the body of the new loop, see [4].

The fact that programs only have a single execution path makes WCET analysis trivial: First, path analysis is superfluous: observing the execution path of any code execution with any input data yields the singleton execution path. Second, the analysis does not need complex and accurate hardware timing models for static WCET analysis. Since programs following this paradigm only have a single path, this singleton path is necessarily the worst-case path. Thus, obtaining the WCET by measurements is possible (either by measurements on the target or on a cycle-accurate hardware simulator) and there is no need to build any specific tools for static analysis. The latter also provides a solution to dealing with new hardware features in the analysis (see above). As the WCET analysis of single-path programs does not require hardware modelling, software developers do not have to wait until tool vendors incorporate the new features into their models in order to perform WCET analysis for their new platforms.

### 4 Conclusion

"Is WCET analysis a non-problem?" is the question posed in the title. To answer this question we investigated whether highly sophisticated WCET analysis techniques are the correct way to deal with the complexity of task timing analysis. We discussed hardware and code design practices that cause complexity and proposed an alternative memory architecture and the single-path programming paradigm as possible ways out.

The answer to the original question seems to be "Yes and No": As long as real-time code is coded for speed rather than temporal predictabiliy and hardware manufacturers continue to use memory hierarchies that rely on speculation then the answer is "no" — and we will certainly have to deal with such systems for at least one more decade. On the other hand, if people become aware of the importance of temporal predictability and build systems correspondingly, then WCET analysis indeed becomes trivial. So the new question to ask is if it will be possible to convince both hardware manufacturers and code designers to change their way of thinking and put temporal predictability first.

#### Acknowledgments

The author would like to thank Raimund Kirner for his valuable comments on an earlier version of the paper and the participants of the 2nd Euromicro WCET workshop for their valuable feedback and discussions.

#### References

- J. Allen, K. Kennedy, C. Porterfield, and J. Warren. Conversion of Control Dependence to Data Dependence. In *Proc. 10th ACM Symposium on Principles of Programming Languages*, pages 177–189, Jan. 1983.
- [2] M. Lee, S. Min, C. Park, Y. Bae, H. Shin, and C. Kim. A Dual-mode Instruction Prefetch Scheme for Improved Worst Case and Average Case Program Execution Times. In *Proc. 14th Real-Time Systems Symposium*, pages 98–105, 1993.
- [3] T. Lundqvist and P. Stenström. Timing Anomalies in Dynamically Scheduled Microprocessors. In *Proc.* 20th IEEE Real-Time Systems Symposium, pages 12– 21, Dec. 1999.
- [4] P. Puschner. Transforming Execution-Time Boundable Code into Temporally Precdictable Code. In Proc. IFIP World Computer Congress, Stream on Distributed and Parallel Embedded Systems, Aug. 2002.
- [5] P. Puschner and A. Burns. Guest Editorial: A Review of Worst-Case Execution-Time Analysis. *Real-Time Systems*, 18(2/3):115–127, May 2000.
- [6] P. Puschner and A. Burns. Writing Temporally Predictable Code. In Proc. 7th IEEE International Workshop on Object-Oriented Real-Time Dependable Systems, pages 85–91, Jan. 2002.