



The abc Group

Optimising AspectJ

abc Technical Report No. abc-2004-3

Pavel Avgustinov¹, Aske Simon Christensen², Laurie Hendren³, Sascha Kuzins¹,
Jennifer Lhoták³, Ondřej Lhoták³, Oege de Moor¹, Damien Sereni¹,
Ganesh Sittampalam¹, Julian Tibble¹

¹ Programming Tools Group
Oxford University
United Kingdom

² BRICS
University of Aarhus
Denmark

³ Sable Research Group
McGill University
Montreal, Canada

November 12, 2004

aspectbench.org

Contents

1	INTRODUCTION	3
2	BACKGROUND AND DEFINITIONS	4
3	OPTIMISING AROUND ADVICE	5
3.1	Implementation Issues	5
3.2	General Implementation	6
3.2.1	Closures	6
3.2.2	Inlining	6
3.3	Around weaving in abc	7
3.3.1	The Generic Implementation	7
3.3.2	Context and Advice Formals	7
3.3.3	Local and anonymous classes	8
3.3.4	Special Cases	8
3.4	Empirical Results	9
4	OPTIMISING CFLOW	10
4.1	Intraprocedural Optimisations	11
4.1.1	cf low without Bound Variables	11
4.1.2	Sharing cf low States	11
4.1.3	CSE of cf low State Retrieval	11
4.2	Interprocedural Optimisations	11
4.2.1	Analysis in abc	12
4.2.2	Call Graph	13
4.2.3	Interprocedural cflow analysis	13
4.3	Empirical Results	15
5	OTHER OPTIMISATIONS	17
6	RELATED WORK	18
7	CONCLUSIONS	18

List of Figures

1	Reweaving in the <i>abc</i> backend	12
2	Jedd code implementing Algorithm 1	15

List of Tables

I	Execution Times and Code Size	9
II	Static intra-procedural optimisation counts	16
III	Benchmark running times (seconds)	16
IV	Static inter-procedural optimisation counts	17

Abstract

AspectJ, an aspect-oriented extension of Java, is becoming increasingly popular. However, not much work has been directed at optimising compilers for AspectJ. Optimising AOP languages provides many new and interesting challenges for compiler writers, and this paper identifies and addresses three such challenges.

First, compiling *around* advice efficiently is particularly challenging. We provide a new code generation strategy for *around* advice, which (unlike previous implementations) both avoids the use of excessive inlining and the use of closures. We show it leads to more compact code, and can also improve run-time performance. Second, woven code sometimes includes run-time tests to determine whether advice should execute. One important case is the *cflow* pointcut which uses information about the dynamic calling context. Previous techniques for *cflow* were very costly in terms of both time and space. We present new techniques to minimize or eliminate the overhead of *cflow* using both intra- and inter-procedural analyses. Third, we have addressed the general problem of how to structure an optimising compiler so that traditional analyses can be easily adapted to the AOP setting.

We have implemented all of the techniques in this paper in *abc*, our AspectBench Compiler for AspectJ, and we demonstrate significant speedups with empirical results. Some of our techniques have already been integrated into the production AspectJ compiler, *ajc* 1.2.1.

1 INTRODUCTION

Aspect-oriented programming is a new programming paradigm that is rapidly growing in acceptance, in large part due to the popularity of AspectJ [4], an aspect-oriented extension of Java that is compatible with existing Java programs. Aspect-oriented programming encompasses many different language features, which can be separated into two categories: *static* features, essentially a version of open classes; and *dynamic* features. While the static features incur no substantial performance penalty, that is not necessarily the case for the dynamic features. New optimisation strategies are required to improve the performance of programs written using the dynamic features of AspectJ.

Challenges The execution model for (the dynamic features of) AspectJ is that execution of the main program is monitored for certain programmer-specified events. When these occur, control is passed to *advice*, a special kind of method that contains the extra code to be executed. Advice can be run before or after the event that triggers it, or it can be executed “around” it (instead of it).

The monitoring of the base program to find points where advice should apply can be expensive. Of course, the AspectJ language has been designed so that most of the conditions that cause advice to be triggered can be determined statically, and so do not induce overheads. A nice explanation from the viewpoint of partial evaluation can be found in [17].

However, some features of AspectJ require the insertion of dynamic checks in the base programs, resulting in performance loss. Furthermore, “around” advice, which we describe in more detail in Section 2, is a powerful language feature that is difficult to implement efficiently. It has been shown that the combination of these factors can result in substantial overhead in AspectJ programs [9]. Here, *overhead* should be understood as the additional cost for matching events to advice, not the execution of advice code itself. In addition, we must consider *space overhead*. This is the size of the extra code introduced when aspects are applied to a program.

These problems highlight the need for novel aspect-specific optimisation strategies.

The *ajc* Compiler The reference implementation of AspectJ, originally developed by the language’s designers, is the *ajc* compiler [4], maintained as part of the Eclipse project. The *ajc* compiler has been designed with compilation speed in mind, and supports incremental compilation. These design aims reduce the opportunities for optimisation, and indeed *ajc* performs only a small amount of intraprocedural analysis and no whole-program analysis. This results in substantial overhead in programs produced by *ajc* when using certain features of AspectJ [9].

The AspectBench Compiler The observation that aspects can introduce substantial overheads has prompted us to build a new compiler for the AspectJ language, the AspectBench Compiler (*abc*) [1], which is freely available under the GNU LGPL. While *abc* uses a weaving strategy similar to *ajc* in many ways, by contrast its design was motivated by two goals: to be *extensible* (so that new features can be added to the input language) and provide a powerful *analysis*

and optimisation framework. This second feature, which aims to reduce the overheads of AOP, is the main focus of this paper.

To achieve this, *abc* uses the Soot program analysis and transformation framework [22] as a backend. Soot provides a typed, three-address intermediate representation (called Jimple) and a library of standard program analyses. In addition, it is straightforward to define new analyses in Soot, both intraprocedural and interprocedural. In particular, the Jimple IR is extremely well-suited to program analysis and transformation.

Contributions This paper is the first systematic study of the analysis and optimisation of aspect-oriented programs, in particular for the AspectJ language. We present the following novel contributions:

- A novel implementation strategy for ‘**around** advice’. Around advice is executed instead of the event that triggered it. Our implementation strategy avoids the use of closures in all but some pathological examples, unlike earlier implementations, which relied heavily on closures. Furthermore, our implementation makes judicious choices about inlining, thus producing compact code.
- A number of intraprocedural optimisations to reduce the overheads of ‘**cflow** pointcuts’. This feature of AspectJ is used to intercept method calls in the dynamic scope of others, and its naive implementation can be very costly.
- An interprocedural analysis to completely eliminate the overheads of using **cflow** in most cases. This analysis and the associated transformations illustrate our final contribution:
- A general technique for leveraging analyses and transformations for pure Java on AspectJ programs. The technique consists of first compiling the program naively, possibly inserting too many dynamic checks (for the applicability of advice) into the base program. We then analyse the resulting intermediate code, and reconsider the decisions to insert dynamic checks, based on the analysis results.

All these contributions have been implemented in *abc* and, at our suggestion, some have also been adopted by *ajc*. We present experiments to confirm that the above techniques can result in dramatic improvements in run time as well as code size. While this paper illustrates the techniques in the context of the AspectJ language, they equally apply to other aspect-oriented programming languages.

Outline This paper is structured as follows. We first briefly introduce the relevant features and terminology of AspectJ in Section 2. Then, we describe the optimisations used in *abc*’s implementation of **around** advice in Section 3. Subsequently, optimisations specific to the **cflow** construct of AspectJ are given in Section 4. Section 5 summarises further efficiency improvements achieved by *abc*. Finally, we give related work in Section 6 and conclude in Section 7.

2 BACKGROUND AND DEFINITIONS

The execution events in an AspectJ program that can be monitored and associated with advice code are called *join points*. A join point is a span of time in the execution of the program. Example join points are a call of a method, an execution of a method body, a field read or an execution of an exception handler. For any particular join point, the textual part of the program executed during the time span of that join point is called the *shadow* of the join point.

AspectJ contains a query language for picking out join points. Such a query is called a *pointcut*. An *advice declaration* consists of an advice kind (**before**, **after**, **around**), a pointcut, and a body of code, called *advice*. The advice is to be executed before, after, or instead of any join point which matches the pointcut. Whenever multiple pieces of advice apply at the same join point, precedence rules determine the order in which they execute.

A pointcut consists of a number of *pointcut designators* connected by boolean connectives. Each pointcut designator is either static (defining a set of join point shadows) or dynamic (defining a run-time condition). Some examples of static pointcuts are **call** and **execution**, which match all calls to and executions of a method matching a pattern; **within**, which matches all join points within class matching a pattern; and **handler**, which matches exception handlers. Some examples of dynamic pointcuts are **args** (*a*), which matches when actual arguments of a method have specified run-time types; **if** (*e*), which matches when the arbitrary Java expression *e* evaluates to true; and **cflow** (*p*), which matches when a join point is within the dynamic scope of a join point matched by pointcut *p*. In addition

to testing runtime conditions, dynamic pointcut designators may also bind context values to variables to make them available to the advice code and to the code in **if** pointcuts. For instance, the **args** pointcut designator may bind the actual arguments.

AspectJ also contains standalone pointcut query constructs. The **declare error** and **declare warning** constructs take a pointcut consisting of only static pointcut designators along with a message. If the pointcut matches any join point shadow, the message is printed as a compile-time error or warning, respectively.

The process of compiling AspectJ programs is known as **weaving**. The base program and advice declarations are woven together into one program which behaves as if the aspects were monitoring execution of the base program, and invoking the relevant advice.

Weaving is a two-step process. The first step, **pointcut matching**, checks for each possible join point shadow in the program and each advice declaration, whether the pointcut could possibly match a join point at that shadow. If so, it constructs a **dynamic residue** of the runtime checks to be performed at the shadow to determine whether the pointcut actually matches.

The second step is **advice weaving**. At each join point shadow where a piece of advice may apply, code is generated to evaluate the dynamic residue and, if it matches, to bind the context values and invoke the advice.

3 OPTIMISING AROUND ADVICE

Most advice in AspectJ programs adds some code to a join point, either before or after it. **around** advice is unique in that it is executed *instead* of each join point it applies to. It can, however, invoke the original join point at any time by using a **proceed** statement.

The following caching aspect demonstrates the usefulness of this interception mechanism. The aspect intercepts calls to a method named *foo* and only executes a call if the result is not in the cache already:

```
aspect Cache {
    pointcut methodsToCache(Object arg) :
        call(Object foo(Object)) && args(arg);
    Object around(Object arg) : methodsToCache(arg) {
        if (!cacheContains(arg))
            setCache(arg, proceed(arg));

        return getCachedValue(arg);
    }
    boolean cacheContains(Object arg) {...}
    void setCache(Object arg, Object value) {...}
    Object getCachedValue(Object arg) {...}
}
```

Like all forms of advice, **around** advice can have arguments. Each of these arguments has to be bound to an exposed context value by the pointcut expression: in the above example, **args(arg)** binds the argument of *foo* to *arg*. A **proceed** statement looks like a method call with the same number of arguments as the advice, and it executes the original join point with the bound context values set to the arguments of the **proceed** call. The **proceed** call can therefore *change* the values of context values (such as method arguments).

3.1 Implementation Issues

A number of difficulties arise in the implementation of **around** advice (and particularly the **proceed** statement), which we briefly outline before giving possible solutions.

The first problem is that **around** advice can apply in multiple places within the base program. For example, consider some **around** advice with pointcut **execution(void foo(..)) || execution(void bar(..))**. This will apply to the bodies of methods *foo* and *bar*, and any occurrence of **proceed** in the advice body will pass control back to the method that was matched. Furthermore, the *context* of the advice applications (the values of locals used in the advised statements) can certainly differ between applications to *foo* and *bar*. This polymorphic behaviour of **proceed** is the main difficulty in its implementation.

Matters are complicated further by the fact that **proceed** can occur in arbitrary places within the advice body, including local and anonymous classes. An important implication is that **proceed** statements may be executed after the control flow has left the advice body. Any context values needed for the **proceed** call must therefore outlive the execution of the advice in this case.

Finally, aspects are not restricted to observing the base program, and in fact advice can apply to other advice. In particular, a piece of **around** advice can apply to the execution of its *own* body, directly or indirectly. Such *circular advice execution applications* are very rare, and usually pathological and a symptom of an error in the program. It is important to observe that these can occur, however, as we shall have to treat such applications as special cases. Note that the application of **around** advice to any other advice other than itself, or to a statement within its body (but not the whole body) is not considered circular. These other cases are common, but circular applications are rare.

3.2 General Implementation

Any piece of advice, regardless of its kind, is turned into a plain Java method (the *advice method*) both by *abc* and *ajc*. The interest lies in the translation of the **proceed** statement, minding the issues described above. The polymorphic behaviour of **proceed**, coupled with the need to store execution context, motivates the use of *closures* as a default and straightforward implementation strategy.

As a preparatory step to implementing **around**, any shadow that is advised by some **around** advice is lifted into a separate method (the *proceed method*) that can be invoked by **proceed**. Note that this is necessary as shadows need not be entire method bodies.

3.2.1 Closures

We shall now give a brief description of the closure strategy.

This approach works by defining a suitable interface type (or class type) for each advice method. All calls to **proceed** are then translated to calls on this closure interface.

```
public interface Around$1 {
    { public [ret-type] proceed$1([arg-type] arg1, ...); }
[ret-type] adviceMethod$1(Around$1 closure,
    [arg-type] arg1, ...) { ...
    [ret-type] result=closure.proceed$1(arg1", ...);
... }
```

For each advice application, a specialised closure type is defined that implements the interface, and for each advice invocation, an instance of the closure class is passed to the advice method. This closure must then call the **proceed** method.

The major drawback of the closure approach is performance. Each time advice is triggered, a closure object has to be allocated on the heap, which can be a significant overhead.

3.2.2 Inlining

In all cases apart from the case of advice applying to itself (either the whole advice or a statement within it), it is possible to avoid closures by duplicating the advice method for each point in the program where the advice applies. This so called inlining (in *ajc* terminology) eliminates the need for polymorphism as the **proceed** statements in this specialised advice method always invoke the same join point shadow, and the join point context can be transferred using method arguments.

While inlining may be a good optimisation in certain situations, the duplication of the advice method for every advice application can lead to code bloat and thus is unsuitable as the general approach.

The **proceed** statement is implemented in *ajc* using these two strategies, a generic strategy based on closures and an inlining strategy that is used in certain cases. We shall now describe *abc*'s approach.

3.3 Around weaving in abc

We present a novel approach for weaving around advice. Our approach is generic: the same strategy can be employed everywhere, and it does not rely on inlining. The only exception is the pathological case of circular adviceexecution applications, described previously. We will return to that case in Section 3.3.4.

When compared to the dual strategy of *ajc*, *abc*'s around weaver never performs significantly worse and in many cases performs significantly better. Particularly,

- whenever *ajc* resorts to closure creation, *abc* can be expected to produce faster code
- when the advice code is big and applies at many locations, *abc* avoids *ajc*'s code bloat
- with circular adviceexecution applications, *abc* produces fewer closure objects and hence faster code.

We shall now describe our weaving strategy in detail.

3.3.1 The Generic Implementation

Instead of creating a closure class for each join point shadow where around advice could apply, *abc* places the proceed code from all join point shadows of a class in a single static proceed method in the class. Each join point shadow at which around advice applies is replaced with a call to the advice body. Into the advice body, we pass a static class ID to identify the class from which we are calling it, and a shadow ID to identify the join point shadow within the class. To implement the **proceed** call, the advice body uses the static class ID (using a switch statement) to select the class from which it was called, and calls the static proceed method in that class. The static proceed method uses the shadow ID to select the shadow whose proceed code it must execute. Because we keep the code of each join point shadow in the class where it originally occurred, there is no need to generate accessor methods for private members.

```
ret-type adviceMethod$1(Around$1 closure,
    int shadowID, int staticClassID, args) {
    ...
    switch (staticClassID) {
    case 0: closure.proceed$1(shadowID, args); break;
    case 1: ShadowClass.proceed$1(shadowID, args); break;
    ...dispatch to other classes to which the advice applies...
    }
    public class ShadowClass {
    public static ret-type proceed$1(int shadowID, args) {
    switch(shadowID) {
    case 0: ... do what the first shadow did...
    case 1: ... do what the second shadow did...
    ...handle further cases...
    }
    }
    }
    ...
}
```

3.3.2 Context and Advice Formals

In addition to calling the right piece of code we must also ensure that values are passed for free variables used by this code (the context). We describe the implementation of context passing next.

Passing context Context from the join point shadow is needed in two places. First, it is used by the code of the shadow itself. Since we have moved this code out from the original join point shadow to the static proceed method, we must pass the required context into this method. Second, the AspectJ pointcut designators **args**, **this**, and **target** allow values from the context of the join point to be bound to formal parameters of the advice, and used in the advice body. Therefore, these values must be passed from the join point shadow site to the advice method that is called.

The dynamic residue of the pointcut guarding the advice may or may not match at run time. If the residue does not match, the advice is not executed, and the static proceed method is called directly from the shadow, so the context values are passed to it directly. If the residue does match, the advice body is called from the shadow, and it may in

turn call the static `proceed` method (if the advice contains a **proceed** call). Therefore, from the shadow, we must pass to the advice body both the context needed by the advice body itself and the context needed by the `proceed` code, so that the advice body can pass it to the static `proceed` method.

One complication is that one advice method can apply to many different join point shadows with different context values. Therefore, we add a sufficient number of parameters of each type to the method implementing the advice to cover the context at all shadows at which the advice may apply. To keep the number of parameters reasonably small, we only add parameters of the types `Object`, `int`, `float`, `double` and `long`, since context values of any type can be converted to one of these types to be passed into the method.

A second complication is that the context value to be bound to an advice parameter may not be known until run time. Consider the following pointcut:

```
void around(Foo x) : args(x,..) || args(...,x)
```

In this case, `x` may be bound to the first argument of a method (if it has type `Foo`), or the last argument.¹ The dynamic residues at the join point shadow determine which part of the pointcut matches, so at the shadow, at run time, we know whether the `x` in the advice should be bound to the first or last parameter.

A third complication is that the advice may modify the context that is to be passed to the `proceed` code for the shadow. The **proceed** expression in the advice body accepts the same number of arguments as there are parameters to the advice body. In the **proceed** call to the original code from the shadow, these arguments replace the context values that were bound to the corresponding advice parameters when the advice was invoked. In the above example, if we call **proceed** (`null`) from the advice, the code of the shadow must be executed with its first or last argument replaced with `null`, depending on which clause of the pointcut matched at the join point before the advice was executed. Since this binding is only known at the dynamic residue and only at run time, we must communicate it from the residue to the advice method, which then communicates it to the static `proceed` method. To do this, at the dynamic residue, we create a bit vector specifying the bindings, and pass it through to the relevant methods.

3.3.3 Local and anonymous classes

As we have observed previously, **proceed** statements can occur in local and anonymous classes within advice methods, and thus the **proceed** invocation can occur after control flow has left the advice method (and so its context must be stored). Our implementation strategy conveniently extends to this case, as all the necessary context is available in the advice method. We simply need to add new fields to the local or anonymous class to hold the context values and initialise these fields when the classes are instantiated.

3.3.4 Special Cases

In the previous sections we have described the generic implementation of **around** advice in *abc*. As mentioned above, in the pathological case of circular advice execution this cannot be used, and we resort to closures instead. Furthermore, for efficiency *abc* also chooses to inline advice methods in some cases (when the advice is small and does not apply numerous times, to avoid code bloat). Here *inlining* is not taken in the *ajc* sense, but rather in the usual sense of substituting the body of the method (in this case, the advice method) at the point where it is called. We now describe these two special cases in greater detail.

Circular advice applications The execution of advice is a join point itself, so advice can apply to the execution of advice (the entire body of the advice method). These advice-on-advice applications can be expressed as a directed graph structure. When weaving into the execution of a method, the whole body of the method is moved into the corresponding `proceed` method. To simplify the weaving process, a topological sort is performed on the graph structure prior to weaving. This ensures that once an `around` advice method has been woven into, it itself is not applied to any join point shadows anymore.

Obviously, a topological sort fails in the presence of cycles in the graph, and the weaver will encounter situations where the advice method to be woven has already been woven into. This is the only case where we resort to the

¹*ajc* as of version 1.2.1 avoids this complication by issuing a compiler limitation error when encountering multiple binding pointcut primitives for the same advice formal.

Benchmark	Time (s)			Size (instr.)		
	<i>abc</i>	<i>abc</i>	<i>ajc</i>	<i>abc</i>	<i>abc</i>	<i>ajc</i>
		(inline)		(inline)		
sim-nullptr	24.0	21.3	22.0	6440	6811	7327
sim-nullptr-rec	24.0	21.5	126.6	6512	14502	9206
weka-nullptr	17.0	15.4	14.1	80199	78949	85590
weka-nullptr-rec	16.9	15.5	43.4	80296	150298	112372
ants-delayed	17.6	17.7	18.3	3768	4361	3785
ants-profiler	22.0	20.2	21.7	7619	17815	13401

Table I: Execution Times and Code Size

creation of closure objects. The semantics of AspectJ dictate that in a cyclical graph, the order of the execution of advice methods is determined by the dynamic residues before any advice is executed. Because those residues and the resulting execution order can be arbitrarily complex, we decided that closures offer a clean, general solution.

As observed previously, cycles in the advice-on-advice application graph are very rare and usually pathological. It therefore seems unnecessary to try to avoid closures under these circumstances. We have, however, strived to minimise the cost of closures; we create specialised closure classes with fields matching the types of context values, whereas *ajc* uses an expensive object array to store context values (requiring boxing of primitive types).

Inlining as an optimisation pass For very small advice, the most efficient strategy can be to inline the advice directly into the join point shadow. This is implemented in *abc* as an optimisation pass running after the around weaver.

The inlining process is implemented as a series of plain Java optimisations. The advice method is first inlined into the join point shadow as a normal Java method. A constant propagator and switch statement folder are then used to remove checks on the static class ID. Finally, the proceed method is inlined also, and its switch statement removed. Since multiple pieces of advice can apply at the same shadow, this whole process must be repeated until there are no calls left to inline.

3.4 Empirical Results

To compare our strategy to *ajc*'s and to experiment with the different tradeoffs of inlining strategies, we experimented with three base programs and three aspects. The base programs are *ants*, an aspect-oriented simulation of an ants colony (following the specification of the ICFP 2004 programming contest) written by one of the authors (OdM) for use in an undergraduate course; *sim*, a discrete event simulator for certificate revocation simulation [2]; and *weka*, part of the weka machine learning library [23]. All benchmarks were run on a dual AMD Athlon 2000+ with 2GB RAM and the Sun J2SE 1.4.2 JVM.

Table I shows the both the execution time (in seconds) and woven code size (in bytecode instructions). For each benchmark we give results for: *abc*, using *abc*'s generic around weaver; *abc* (inline), the same as *abc*, but with the postpass inlining optimization; and *ajc*, the result given by *ajc*'s around weaver (which is either closure-based or inlining, depending on the benchmark). For each benchmark, we have put the fastest time and the smallest code size in bold. Note that in almost all cases either *abc* or *abc* (inline) gives the fastest code and that *abc*'s code size is smallest, sometimes by a significant margin.

To compare *abc*'s generic weaving strategy to *ajc*'s closure-based and inlining strategies, we applied two versions of the *nullptr* aspect [3] to two base programs, *sim* and *weka* (the first four lines in Table I). The *nullptr* aspect is a very simple **around** aspect for enforcing coding standards that we found on the web when searching for examples of aspects. It simply checks for methods returning and issues error messages in the cases where null is returned. We used two different versions of the aspect, a recursive one (the original form) where the advice applies to itself and a non-recursive version where we explicitly use *!within(...)* to avoid matches within the body of the advice. The *ajc* compiler uses closure object creation for the first case (because of the recursion) and inlining for the second case, whereas *abc* uses its generic implementation for both cases.

Comparing the execution times for the non-recursive (*sim-nullptr* and *weka-nullptr*) versions to the times for the recursive versions (*sim-nullptr-rec*) and (*weka-nullptr-rec*), we can see that the execution time and code size for *abc*

is almost the same, whereas *ajc* produces much slower code for the recursive versions (almost 6 times slower for *sim* and 3 times slower for *weka*). For the the non recursive cases, *abc* is slightly slower than *ajc*, but *abc* (inline) provides further performance improvement at the expense of code size.² Even with our inliner the *ajc* inlining strategy is slightly faster for the *weka-nullptr* case.

From these experiments we can see that *abc* is fairly insensitive to whether the advice is recursive or not, but *ajc* pays a huge penalty when it must switch to an explicit closure strategy. *abc*'s behaviour is beneficial since programmers often make their advice recursive by accident and they need not pay a performance penalty. Furthermore, *abc* allows for further performance improvement with the specialized postpass inliner.

There are other situations where the *ajc* weaver uses closures, which is demonstrated by the *ants-delayed* benchmark. This benchmark uses the *DelayOutput* aspect which captures calls to output methods and delays these calls until the end of the base program. This is accomplished using a local class of type *Runnable* in the advice method that calls **proceed** in its *run()* method. The *ajc* weaver has to instantiate closure objects in addition to the instances of the local class. Our weaving strategy avoids this, which explains why the *abc* results are slightly faster.

To demonstrate the adverse effects of a naive inlining strategy, we applied a profiling aspect to our ants base program (*ants-profiler*). The profiling aspect contains a relatively big piece of around advice that is applied to the execution of every method in the base program. Note that *ajc*'s inlining strategy doubles the size of the resulting class files due to the inlining of the advice code into every method. Note that this increase in code size can also be observed with *abc*'s inlining strategy. However, with our weaving strategy, inlining is an entirely optional optimisation and turning it off only has a slight effect on the efficiency of the resulting program. In *ajc*'s case, the only alternative to inlining is the use of closures with the dramatic effects on performance shown in the table. Furthermore, with our approach we can selectively inline and we are actively working on inlining heuristics specific to the around weaver.

4 OPTIMISING CFLOW

The **cfow** pointcut picks out join points that fall within the dynamic scope of certain events. Specifically, for any pointcut *p*, **cfow**(*p*) applies at a point in the execution of the program if *p* matches *some* state in the call stack at that program point. If *p* contains variables to be bound, then these are bound to the actual values found in the match nearest the top of the call stack. For example, the pointcut

$$\text{call}(*foo()) \&\& \text{cfow}(\text{call}(*bar(*)) \&\& \text{args}(x))$$

matches all calls to *foo* that occur within the dynamic scope of a call to *bar*, and binds *x* to the value of the argument of the last call to *bar*.

It is clear that the use of **cfow** pointcuts requires, in general, the insertion of dynamic tests in the program to test the current state against the condition **cfow**(*p*). The naive implementation of **cfow** associates a state with each **cfow** pointcut (this implementation is described in [17]) and updates this incrementally. The state is a stack of variable bindings that represents an abstraction of the call stack. Each time a join point that matches *p* is entered, a new item is pushed onto the stack, with all the variables in *p* bound to the appropriate values. When this join point is left, the top of the stack is popped. Finally, to check whether **cfow**(*p*) applies at a program point, it suffices to check whether or not the stack is empty; if it is non-empty then the pointcut applies and the appropriate variable bindings can be found on top of the stack.

The implementation of **cfow** (as described above and used in *ajc*) is clearly expensive, both because of the need to update the state (which happens every time *p* applies) and because of the dynamic tests inserted (which can, in the worst case, apply everywhere). Performance experiments confirm that the overhead introduced is substantial [9].

We introduce a number of optimisations for **cfow**, all implemented in *abc*. We first show a number of simple, intraprocedural optimisations that reduce the overhead substantially. Then, we show how the overhead can be entirely eliminated in many common cases by an interprocedural analysis. Finally, we give empirical measurements showing that the optimisations are very effective.

²In some cases the inlined code can be slightly smaller because it removes calls that required a large number of arguments.

4.1 Intraprocedural Optimisations

The simple **cflow** optimisations focus on eliminating the more obvious inefficiencies in updating the state and checking for applicability. They are straightforward but quite effective.

4.1.1 cflow without Bound Variables

The first optimisation applies to pointcuts of the form **cflow**(p), where p does not bind any values. In this case, the state of the **cflow**(p) reduces to a stack of empty sets of variable bindings. In *ajc* 1.2, this is represented by a stack of arrays of length 0.

We improve on this in the obvious way, by replacing the stack with an integer counter that is incremented and decremented when p is entered and left respectively. This avoids repeated allocations of empty arrays. The case of a parameterless **cflow** appears to be quite common, so this optimisation is widely applicable.

4.1.2 Sharing cflow States

Another optimisation that *abc* performs is to *share* the state update and query code between related (or identical) **cflow** pointcuts whenever possible. Consider the following pointcuts:

```
call(* bar()) && cflow(call(* foo(..) && args(t, *, *))  
call(* bar()) && cflow(call(* foo(..) && args(*, s, *))
```

A naive implementation would keep a stack for each **cflow** pointcut, and update and query them independently. We optimise this by observing that a single **cflow** pointcut can be written that covers the two existing instances. In this case, it is **cflow**($call(* foo(..) \&\& args(l_1, l_2, *)$) (where l_1 and l_2 are fresh variables). Note that this binds variables used in either one of the **cflow** pointcuts in the original program.

The implementation of **cflow** in *abc* attempts to *unify* each pair of **cflow** pointcuts that it finds. Unification of two pointcuts succeeds if the pointcuts are syntactically equivalent with the exception of free variables, and returns a pointcut that carries enough state to cover both pointcuts (as in the above example).

In general, this sharing of state can improve performance substantially. In fact, cases similar to the above arise frequently due to inlining of named pointcuts, a strategy used both in *ajc* and *abc*. An added benefit is that some method bodies can become smaller when this is performed (by avoiding duplication of bookkeeping code). We present empirical measurements of the performance improvements in Section 4.3.

4.1.3 CSE of cflow State Retrieval

The final simple **cflow** optimisation is the caching of **cflow** state objects (stacks or counters). The state of a given pointcut **cflow**(p) is thread-local, as it is an abstraction of the call stack. Multiple copies are therefore kept, one for each thread, and any operation on **cflow** state (updating or checking) involves retrieving the copy valid for the current thread (in the worst case, a hash table lookup).

In general, multiple operations on the same **cflow** state can occur within the body of the same method. In fact, updates to the state of a **cflow** are always paired (the state is updated when entering and leaving a join point), so in most cases the state is retrieved at least twice in any method in which it is needed at all.

We can therefore improve on the original implementation by retrieving the appropriate state object only when it is first used in a given method, and storing it for future uses in the same method.

4.2 Interprocedural Optimisations

The optimisations that we have described above reduce the overhead associated with **cflow**, but this can still be substantial. Since the **cflow** construct depends on dynamic properties of the program in general, it is impossible to eliminate such overhead entirely. However, many uses of **cflow** can be statically determined, at least at some program points. To take a simple concrete example, the pointcut **cflow**($call(* foo())$) matches all points in the execution of the program

within the dynamic scope of a call to *foo*. It is possible to determine statically that some program points can *never* be in the dynamic scope of *foo*, and that some program points *always* execute in its scope. At each such program point the **cf**low pointcut is statically known to be true or false, so the dynamic check can be eliminated. In addition, eliminating such dynamic matching code can allow the compiler to eliminate some of the state-updating code for this **cf**low (if its effects can no longer be observed after dynamic checks are removed).

Our empirical results in Section 4.3 show that **cf**low pointcuts can indeed be statically determined in almost all cases. We will now describe the analysis used in *abc* to achieve this.

The idea for this analysis was introduced in [18] for a simple procedural language. The analysis has been adapted to the much wider context of AspectJ and implemented within *abc*. It requires an interprocedural analysis, but has two substantial advantages:

- It eliminates the overhead for **cf**low completely in many common cases, and
- In those cases, it allows **cf**low to be used in constructs that require static pointcuts (such as **declare warning**).

4.2.1 Analysis in *abc*

One of the design goals of *abc* was to make it possible to analyze the code being woven, and use the analysis results to optimise the weaving process to produce more efficient code. In particular, we wanted to be able to leverage the many analyses existing for Java code, without having to rewrite all of them to be specific to AspectJ. Therefore, *abc* includes a hook to perform analyses on the Jimple code produced immediately after weaving, optimise the naive weaving instructions originally produced by the matcher, and then repeat the weaving process on the original code using the optimised weaving instructions. Because the woven code being analyzed has no AspectJ-specific constructs, it is possible to apply standard analyses already in Soot. Of course, we also implement analyses and optimisations specific to AspectJ, but these are greatly simplified by being able to use the results of Java analyses.

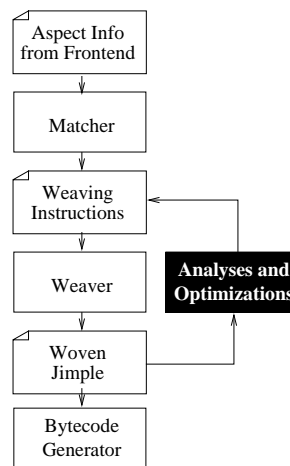


Figure 1: Reweaving in the *abc* backend

The structure of the *abc* backend which makes these analyses and optimisations possible is shown in Figure 1. In normal operation, the phases are executed from top to bottom. The matcher takes as input information about all the pointcuts and advice from the front-end, and produces a set of weaving instructions which specify where in the code they should be woven. The weaver executes these instructions, producing woven Jimple, which is finally translated to bytecode. When we want to perform optimisations, we feed the woven Jimple back into our analyses, and use their results to optimise the weaving instructions that were used to produce the woven code. The weaver can then execute the optimised weaving instructions on the original Jimple code. Note that this mechanism requires saving a copy of the original Jimple code prior to the original weaving pass, which *abc* does.

So far, we have implemented an interprocedural **cf**low analysis and a `thisJoinPoint` escape analysis, but the approach is general; other analyses can be added to the box labelled Analyses and Optimisations in Figure 1.

It is possible that an analysis will produce more precise results if executed not on the naively woven Jimple, but on Jimple woven using optimised weaving instructions produced by an earlier pass of the analysis. Therefore, *abc* allows

the Analyses and Optimisations feedback loop to be repeated if desired.

4.2.2 Call Graph

Estimating which shadows may or must be in the cflow of other shadows requires a call graph approximating which methods may be called from which call sites. We base our analyses on a conservative call graph: every method invocation possible at run-time must be included in the call graph. Call graph construction for object-oriented languages like Java has been the subject of a significant amount of research (*e.g.* [5, 8, 20, 21]). Rather than reinvent the wheel, we construct call graphs using Paddle, a successor of Spark [12, 13], the points-to analysis and call graph construction framework available in Soot [22].

In Java, most method calls are virtual, meaning that the method invoked depends on the run-time type of the receiver object. The treatment of virtual calls is one of the key features distinguishing different call graph construction algorithms. The Paddle framework allows us to experiment with call graphs constructed using algorithms ranging from CHA [8], which conservatively assumes that receivers could have any type admitted by their declared type, to using a subset-based points-to analysis to compute possible run-time receiver types.

Some applications of call graphs, such as devirtualisation, only require call edges for explicit invoke instructions present in the code. However, because methods invoked implicitly by the VM are defined to be in the cflow of their calling context, our call graph must include these implicit calls. In particular, we include implicit calls to static initialisers [16, section 2.17.4], calls through the `PrivilegedAction` interface, and implicit constructor calls by the `Class.newInstance` method. For the latter, the user provides a list of all classes that may be instantiated using reflection. To ensure that this list (and our call graph) is complete, we insert code into methods not reachable in the call graph to abort execution and alert us to the error. Paddle handles these tricky but important details for us; we do not need to consider them explicitly in our cflow analysis.

One type of implicit method invocation which we specifically exclude from the call graph used for cflow analysis is the invocation of the `run` method of newly created threads. In *abc*, we strive to be consistent with the AspectJ language as specified by the *ajc* implementation. In *ajc*, **cfow** stacks are maintained separately for each thread, so the code executed by a thread is not considered to be in the cflow of the code that created the thread. The call graph produced by Paddle includes a type for each call edge, so we can remove these unwanted invocation edges.

4.2.3 Interprocedural cflow analysis

Desired optimisation The customary implementation of a **cfow** pointcut expression **cfow**(*p*) incurs overhead at two kinds of shadows. First, at each shadow matching *p*, a cflow stack is pushed and popped to indicate when we are in the dynamic scope of the **cfow**. We denote these shadows with the term *update shadow*. Second, at each shadow where the **cfow** pointcut could possibly match, we insert a dynamic residue to test whether the **cfow** stack is non-empty. We denote these shadows with the term *query shadow*.

We wish to perform two kinds of optimisation. First, if we can determine **cfow** stack emptiness at a query shadow statically, we can remove the dynamic residue at the query shadow, and possibly other code that becomes unreachable. Second, if we can prove that a **cfow** stack update operation will not be observed by a stack query within the dynamic scope of an update shadow, we can remove the stack update operations at the update shadow.

Analysis information required For each update shadow *sh* in the program, we define two sets of instructions to be computed, **mayCflow**(*sh*) and **mustCflow**(*sh*). **mayCflow**(*sh*) contains every instruction *i* in the program such that when *i* is executed, we may be in the dynamic scope of *sh*. That is, *i* may execute after the push operation of *sh* has been performed, but before the corresponding pop operation has been performed. **mustCflow**(*sh*) contains every instruction *i* such that whenever *i* is executed, we must be in the dynamic scope of *sh*.

Whenever a query shadow is not in **mayCflow**(*sh*), we replace the dynamic test with a constant false pointcut expression.³ Any query shadow in **mustCflow**(*sh*) is replaced with a constant true pointcut expression.

³The **cfow** expression may be part of a more complicated pointcut expression. Constant folding of pointcut expressions is done in a separate phase prior to weaving.

In addition, we calculate a subset **necessaryShadows** of update shadows whose effect may be observed at a query shadow. Each update shadow $sh \in \mathbf{necessaryShadows}$ satisfies two properties. First, some query shadow qsh that has not been resolved statically may occur in the dynamic scope of sh (i.e. $qsh \in \mathbf{mayCflow}(sh)$). Second, sh may occur outside the dynamic scope of all update shadows for the same **cflow** stack (i.e. $\nexists sh'.sh \in \mathbf{mustCflow}(sh')$). This second condition enables us to mark as unnecessary those update shadows at which the stack is always already non-empty.

The optimisations become more complicated when the **cflow** binds arguments because, in this case, each query shadow not only tests whether the stack is non-empty, but also observes the entry at the top of the stack. We can still resolve statically those query shadows not in **mayCflow**(sh), since we know that the stack would always be empty when they are executed. However, at the query shadows where we know the stack is non-empty, we must keep the dynamic residues which read the entry from the stack. In addition, we can no longer remove update shadows just because they are in the **mustCflow** of some other update shadow which will make the stack non-empty, because we also need the correct entry to be pushed onto the stack in addition to the stack being non-empty.

Computing analysis information The exact extent of a **cflow** shadow depends on subtle details of advice precedence and the distinction between **cflow** and **cflowbelow**, and the weaver must respect these details when weaving the **cflow** stack update operations. Because we perform the analysis on the woven code, we need not consider these details; we simply consider each **cflow** shadow to start immediately after the point where the weaver wove the **cflow** push instruction, and end immediately before the corresponding **cflow** pop instruction. We need to unambiguously classify every instruction in the method as being either within or outside the **cflow** shadow. This requires that there be no jumps into or out of the shadow, which would bypass the push or pop instruction.

Due to details of the weaving process, this requirement is always satisfied, except in the case when the argument p of the **cflow** expression $\mathbf{cflow}(p)$ is not entirely static, and requires a dynamic residue. In this case, the weaver generates the dynamic test at the update shadow. If the pointcut p does not match, we do not enter the dynamic scope of the **cflow**, so a conditional jump skips the stack update operations. Therefore, when p is not entirely static, the instructions between the push and pop may execute within or outside the dynamic scope of the **cflow**. Since no instruction can be guaranteed to execute only in the dynamic scope of the **cflow**, $\mathbf{mustCflow}(sh)$ is the empty set in this case.

Algorithm 1 is used to compute $\mathbf{mayCflow}(sh)$. It begins with the statements in the intra-procedural shadow of sh . Then, it adds the statements of all methods that may be called from a statement already in the set, until a fixed point is reached.

Algorithm 1 compute $\mathbf{mayCflow}(sh)$

```

mayCflow  $\leftarrow \{st \mid st \text{ is in intra-procedural shadow of } sh\}$ 
repeat
  for all methods  $m \mid \exists st \in \mathbf{mayCflow}.st \text{ may call } m$  do
    mayCflow  $\leftarrow \mathbf{mayCflow} \cup \text{set of statements in } m$ 
  end for
until mayCflow does not change

```

We have implemented all of the inter-procedural **cflow** analyses using Jedd [14], an extension of Java for expressing analyses using binary decision diagrams (BDDs), which it abstracts as relations. We chose to implement the analyses in Jedd for two reasons. First, they can be expressed in Jedd concisely and clearly. As an example, Figure 2 shows the Jedd implementation of Algorithm 1. Notice that the implementation closely mirrors the algorithm. Second, although the sets computed in the analyses may become quite large, they are likely to share many similarities. BDDs make it possible to represent these large sets compactly.

The set $\mathbf{mustCflow}(sh)$ is computed using Algorithm 2. If sh has a dynamic residue, it must return the empty set. Otherwise, it starts with all statements in the program, and removes statements that can be reached from the entry points of the call graph without passing through sh . The statements to be removed are computed by starting with the entry points, and adding statements of methods called from the set computed so far, but excluding statements in the intra-procedural shadow of sh , until a fixed-point is reached.

Computation of **necessaryShadows** is shown in Algorithm 3. We begin with all the query shadows, and remove

```

<stmt> mayCflow(Shadow sh) {
  <stmt> mayCflow = stmtsWithin(sh);
  <stmt> old;
  do {
    old = mayCflow;
    <method> targets =
      mayCflow{stmt} <> callTargets{stmt};
    mayCflow |=
      targets{method} <> stmtsIn{method};
  } while( mayCflow != old );
  return mayCflow;
}

```

Figure 2: Jedd code implementing Algorithm 1

Algorithm 2 compute **mustCflow**(*sh*)

```

if sh has a dynamic residue then
  mustCflow  $\leftarrow \emptyset$ 
else
  mustCflow  $\leftarrow$  set of all statements
  shadowStmts  $\leftarrow \{st \mid st \text{ is in intraprocedural shadow of } sh\}$ 
  targets  $\leftarrow$  set of entry points of call graph
  repeat
    targetStmts  $\leftarrow$ 
       $\{st \mid \exists m \in \text{targets}. st \text{ is an statement in } m\} \setminus \text{shadowStmts}$ 
    mustCflow  $\leftarrow$  mustCflow  $\setminus$  targetStmts
    targets  $\leftarrow \{m \mid \exists st \in \text{targetStmts}. st \text{ may call } m\}$ 
  until mustCflow does not change
end if

```

those known statically to be false. Unless the **cflow** binds arguments, we can also remove those known statically to be true. This leaves us with the query shadows that will be tested dynamically. The necessary shadows are now those update shadows in whose **mayCflow** any dynamic query shadow appears. Unless the **cflow** binds arguments, we can also remove those update shadows which are already in the **mustCflow** of another update shadow.

Algorithm 3 compute **necessaryShadows**

```

queries  $\leftarrow$  set of all query shadows  $\cap \bigcup_{sh} \text{mayCflow}(sh)$ 
if cflow does not bind arguments then
  queries  $\leftarrow$  queries  $\setminus \bigcup_{sh} \text{mustCflow}(sh)$ 
end if
necessaryShadows  $\leftarrow \{sh \mid \exists st \in \text{queries}. st \in \text{mayCflow}(sh)\}$ 
if cflow does not bind arguments then
  necessaryShadows  $\leftarrow$ 
    necessaryShadows  $\setminus \bigcup_{sh} \text{mustCflow}(sh)$ 
end if

```

4.3 Empirical Results

The **cflow** optimisations we present in this paper have been empirically validated in two different AspectJ compilers. First, we have implemented all the optimisations in our *abc* compiler. Second, we suggested them to the *ajc* team, and they have implemented counters (Section 4.1.1) and sharing (Section 4.1.2) in *ajc* 1.2.1.

We tested the **cflow** optimisations on benchmarks from a wide range of sources. We list the benchmarks and their sizes (non-comment SLOC) in the first column of Table II. Figure is a demo from the AspectJ documentation. Quicksort is the example from [18] with modifications suggested by Gregor Kiczales. Sablecc is a compiler written using the SableCC compiler generator, with an aspect applied to count memory allocations in each of its phases. The base programs *ants*, *certrevsim(sim)* and *weka* were introduced in Section 3.4. Law of Demeter [15] is a style-checking

aspect that we have applied to two code bases: Certresim and weka. Cona [19] is a framework for specifying and checking pre- and post-conditions using aspects. We applied it to the stack example mentioned in the paper, and to the simulator.

Benchmark	SLOC	No opt.	Intra-proc	
		Stacks	Stacks	Counters
figure	94	5	0	1
quicksort	72	2	0	1
sablecc	31233	2	0	2
ants	939	1	1	0
LoD-sim	1586	13	0	1
LoD-weka	3912	13	0	1
Cona-stack	291	10	0	1
Cona-sim	1942	46	0	8

Table II: Static intra-procedural optimisation counts

In Table II, we present the static effects of our intra-procedural optimisations implemented in *abc*. The column labelled “no opt. stacks” shows the number of different stacks before our optimisations; the “intra-proc” column shows the number of stacks and counters after intra-procedural optimisations have been applied. In most cases, sharing reduces the number of **cfow** stacks (or counters) significantly, often down to one. In all benchmarks except ants, all **cfow** stacks are replaced with counters. A counter cannot be used for ants because the **cfow** pointcut binds a value.

Benchmark	<i>abc</i>			<i>ajc</i>	
	No opt.	Intra-proc	Inter-proc	1.2	1.2.1
figure	2139.32	39.53	2.04	506.28	219.71
quicksort	124.88	27.08	27.21	124.07	29.26
sablecc	34.23	24.01	21.16	30.42	24.87
ants	34.40	32.70	13.17	34.66	34.30
LoD-sim	615.31	28.32	22.39	4801.64	39.05
LoD-weka	1969.77	83.35	62.64	2372.80	126.22
Cona-stack	1153.66	32.93	23.51	1151.63	65.25
Cona-sim	77.12	73.10	71.28	76.77	69.76

Table III: Benchmark running times (seconds)

We present the benchmark running times in Table III. The middle section lists the running times of benchmarks compiled using *abc* with **cfow** optimisations disabled, with the intra-procedural optimisations described in Section 4.1, and with the inter-procedural optimisations described in Section 4.2. The rightmost section lists running times when the benchmarks are compiled with *ajc* versions 1.2 and 1.2.1. Between these two releases, two of the intra-procedural optimisations presented in this paper, counters (Section 4.1.1) and sharing (Section 4.1.2), were added to *ajc* in response to our suggestions.

Using the *abc* compiler, the speedups due to our intra-procedural optimisations are very significant (up to 54-fold) not only for small benchmarks (*e.g.* figure, quicksort), but also for large benchmarks which use **cfow** (*e.g.* the LoD benchmarks). We observe similar speedups with the *ajc* compiler between version 1.2 and 1.2.1, in which intra-procedural **cfow** optimisations were added.

Static results of our inter-procedural **cfow** analysis are shown in Table IV. The “query shadows” column shows, for each **cfow** pointcut designator (corresponding to a stack or counter), the total number of query shadows and, of those, how many the analysis determined to be unreachable, how many are determined to never or always match, and how many cannot be determined statically and therefore still require a dynamic test. The “update shadows” column shows the total number of update shadows and the number that the analysis determines to be necessary, and must remain as dynamic updates even after the analysis.

With the exception of one **cfow** pointcut designator in sablecc, the analysis was able to statically determine the outcome of all **cfow** queries, and therefore entirely remove the dynamic updates and queries of the **cfow** stacks or counters. The imprecision in the sablecc case is due to query shadows in a static initialiser; to deal with this case, we are developing a simple analysis to reduce the number of spurious static initialiser edges in our call graph.

Even though the **cfow** pointcut in ants binds an argument, we can eliminate it because it is never queried. This

Benchmark	Query shadows					Update shs.	
	Total	Unreach.	Never	Always	Dynamic	Total	Dynamic
figure	6	0	2	4	0	6	0
quicksort	6	0	2	4	0	3	0
sablecc	985	388	299	298	0	698	0
	985	388	332	260	5	1	1
ants	84	0	84	0	0	1	0
LoD-sim	1313	798	515	0	0	41	0
LoD-weka	7031	3501	3530	0	0	41	0
Cona-stack	16	0	14	2	0	27	0
Cona-sim	2	0	2	0	0	2	0
	3	3	0	0	0	18	0
	4	3	1	0	0	31	0
	0	0	0	0	0	2	0
	7	5	2	0	0	20	0
	0	0	0	0	0	6	0
	4	0	4	0	0	5	0
	0	0	0	0	0	3	0

Table IV: Static inter-procedural optimisation counts

is because the pointcut is being used as an assertion to find an error condition. By determining that the **cflow** never matches, we have statically verified the assertion. The success of the static analysis inspired us to begin experimenting with AspectJ extensions to allow “dynamic” pointcuts such as **cflow** in “static” **declare error** constructs. This provides a way for a programmer to specify properties of the program to be checked. When the analysis cannot prove the properties at compile time, a warning is issued and a run-time check inserted.

We were pleasantly surprised that the inter-procedural analysis was so effective in resolving **cflow** statically. To ensure that these analysis results are indeed correct, we ran all the benchmarks with a special dynamic residue woven in to check that the static analysis results always agreed with the run-time behaviour.

The performance improvements due to the removal of update shadows, query shadows, and the code of unreachable advice are shown in the “*abc* inter-proc” column of Table III. On benchmarks making significant use of **cflow**, both small (*e.g.* figure) and large (*e.g.* LoD), these optimisations provide large speedups, even on top of the already large speedups from the intra-procedural optimisations and the use of cheap **cflow** counters. Furthermore, when the **cflow** binds an argument, the cheap counters cannot be used, so the inter-procedural optimisations enable the removal of expensive **cflow** stacks, resulting in a 2.5-fold speedup in ants.

5 OTHER OPTIMISATIONS

abc implements a number of other optimisations, and we now give a brief overview of these. They fall into two groups: those that affect reflective access to the current join point in pointcuts and advice; and the backend optimisations afforded by the Soot framework.

It is very common for advice to only access the static part of the *thisJoinPoint* object. *abc* performs a simple, conservative analysis to determine whether this is the case, and if so, it replaces the use of *thisJoinPoint* by a smaller structure, *thisJoinPointStaticPart*, which can be completely computed at compile-time. *thisJoinPointStaticPart* is in fact part of the AspectJ language; this optimisation thus relieves the programmer from the burden of deciding whether *thisJoinPoint* or *thisJoinPointStaticPart* is preferable. In cases where we do need to keep the dynamic version, it is initialised lazily. In particular it is not constructed prior to the dynamic pointcut matching, as such construction might turn out to be in vain if the pointcut fails to match. These optimisations are also present in *ajc*.

After all the aspect-specific transformations are complete, *abc* runs a number of generic optimisations for Java bytecode, which are part of the Soot framework. These are: local packer (‘register allocation’ on bytecode), copy and constant propagation, common subexpression elimination, partial redundancy elimination, dead assignment elimination unreachable code elimination, branch simplification. We have found these (in particular the local packer) to be effective in producing good code for aspects.

6 RELATED WORK

This work is the first general study of analysis and optimisation strategies for aspect-oriented languages in general, and for AspectJ in particular. As a consequence, the amount of related work is rather sparse. There are however a number of other industrial strength implementations of aspect-orientation, and we discuss these here.

The *ajc* Compiler The reference implementation of AspectJ (and in fact the only other implementation of the language) is the *ajc* compiler. The weaving strategies of *ajc* and *abc* are similar, except for the optimisations described in this paper. Following the early success of our optimisations in *abc*, two of them (**cflow** counters and sharing of **cflow** stacks) have been incorporated into *ajc* 1.2.1. Further details on the implementation of weaving in *ajc* (similar to weaving in *abc* except for the optimisations described here) can be found in [10].

Other AOP Systems There are a number of other systems besides AspectJ that support the use of aspects. Perhaps the most successful of these is AspectWerkz [7]; its features are in fact very similar to those of AspectJ. Aspects are however deployed using annotations or scripts, rather than in an extension of the Java language. Unlike AspectJ, AspectWerkz supports dynamic weaving: enabling new aspects at runtime, and also disabling them:

The AspectWerkz system may however also be used in off-line mode, in the same way as *ajc* or *abc*. Because of its focus on runtime weaving, AspectWerkz employs an event-based implementation of join points, where advice can register as a listener. In preliminary experiments, we have found this strategy leads to a slow-down of a factor of 9 or more compared to *ajc* or *abc*. Because of this huge gap and the different aims of dynamic weaving, we focused on the most popular static weaving system, which is *ajc*. The other leading AOP system is JBoss [11], and this employs a similar implementation strategy to AspectWerkz. It is evident that both these systems could benefit by the optimisations presented here, when used in off-line mode. For weaving at runtime, it would appear that our intraprocedural optimisations may be helpful. The same applies to efforts to support aspects in a modified JVM, as in [6].

7 CONCLUSIONS

The field of optimising compilers for AOP languages is just starting, but we believe that this area will provide many interesting problems and challenges that can be met with both existing and new compiler optimisation technology.

In this paper we have presented three main contributions to the field in the context of a new optimising compiler for AspectJ, *abc*. We have designed and implemented a new strategy for weaving **around** advice which aims to avoid both the code size explosion of a pure inlining approach and the time and space overhead of an explicit closure-based approach. Our experimental results demonstrate that this technique works very well, it is much more efficient than the closure-based approach, and produces much less code than the inlining-based approach.

Our second major contribution was to show how to reduce or eliminate the large overheads associated with **cflow**. We gave some intra-procedural techniques that are both relatively simple and very effective at reducing large overheads for the common case. These optimisations have already been adopted by the implementors of the *ajc* compiler. We then showed that we can go even further by applying inter-procedural analyses that can statically approximate dynamic **cflow** properties. Our experimental results show that in many cases we can completely eliminate the **cflow** overhead.

Finally, the implementation strategies presented here showcase a novel methodology for defining new program analyses and efficiency-improving program transformations for aspect-oriented languages. In particular, the inter-procedural **cflow** analysis shows that *reweaving* is a useful technique. In reweaving, aspects are woven first naively into the base program, the resulting program is analysed and the results of the analysis are used to guide subsequent weaving phases (so that better code can be produced). In general, reweaving can be iterated multiple times. The *abc* compiler was designed specifically to support reweaving and thus can serve as a workbench for developing new optimising transformations of AspectJ.

Acknowledgments

This work was supported, in part, by NSERC, EPSRC and IBM.

References

- [1] abc. The AspectBench Compiler. Home page with downloads, FAQ, documentation, support mailing lists, and bug database. Available from URL: <http://aspectbench.org>.
- [2] André Årnes. PKI certificate revocation. Available at <http://www.pvv.ntnu.no/~andream/certrev/>.
- [3] R. Dale Asberry. Aspect oriented programming (AOP): Using AspectJ to implement and enforce coding standards. <http://www.daleasberry.com/newsletters/200210/20021002.shtml>, 2002.
- [4] Eclipse AspectJ. The AspectJ Eclipse Project. <http://eclipse.org/aspectj>.
- [5] David F. Bacon and Peter F. Sweeney. Fast static analysis of C++ virtual function calls. In *OOPSLA 1996*, pages 324–341, 1996.
- [6] Christoph Bockish, Michael Haupt, Mira Mezini, and Klaus Ostermann. Virtual machine support for dynamic join points. In *AOSD 2004*, 2004.
- [7] Jonas Bonér. Aspectwerkz — dynamic AOP for java. Available from URL: http://codehaus.org/~jboner/papers/aosd2004_aspectwerkz.pdf, 2004.
- [8] Jeffrey Dean, David Grove, and Craig Chambers. Optimization of object-oriented programs using static class hierarchy analysis. In *ECOOP '95*, volume 952 of *LNCS*, pages 77–101, 1995.
- [9] Bruno Dufour, Christopher Goard, Laurie Hendren, Oege de Moor, Ganesh Sittampalam, and Clark Verbrugge. Measuring the dynamic behaviour of AspectJ programs. In *OOPSLA 2004*, October 2004.
- [10] Erik Hilsdale and Jim Hugunin. Advice weaving in AspectJ. In K. Lieberherr, editor, *AOSD 2004*, 2004.
- [11] jboss. JBoss Aspect Oriented Programming. Home page with downloads, documentation, wiki. <http://www.jboss.org/index.html?module=html&op=userdisplay&id=developers/projects/jboss/aop>.
- [12] Ondřej Lhoták. Spark: A flexible points-to analysis framework for Java. Master's thesis, McGill University, December 2002.
- [13] Ondřej Lhoták and Laurie Hendren. Scaling Java points-to analysis using Spark. In G. Hedin, editor, *CC 2003*, volume 2622 of *LNCS*, pages 153–169. Springer, April 2003.
- [14] Ondřej Lhoták and Laurie Hendren. Jedd: a BDD-based relational extension of Java. In *PLDI 2004*, pages 158–169, 2004.
- [15] Karl Lieberherr, David H. Lorenz, and Pengcheng Wu. A case for statically executable advice: checking the law of demeter with aspectJ. In *AOSD 2003*, pages 40–49, 2003.
- [16] Tim Lindholm and Frank Yellin. *The Java Virtual Machine Specification*. Addison-Wesley, second edition, 1999.
- [17] Hidehiko Masuhara, Gregor Kiczales, and Chris Dutchyn. A compilation and optimization model for aspect-oriented programs. In *CC 2003*, volume 2622 of *LNCS*, pages 46–60, 2003.
- [18] Damien Sereni and Oege de Moor. Static analysis of aspects. In *AOSD 2003*, pages 30–39, 2003.
- [19] Therapon Skotiniotis and David H. Lorenz. Cona: aspects for contracts and contracts for aspects. In *OOPSLA 2004 Companion*, pages 196–197, 2004.

- [20] Vijay Sundaresan, Laurie Hendren, Chrislain Razafimahefa, Raja Vallée-Rai, Patrick Lam, Etienne Gagnon, and Charles Godin. Practical virtual method call resolution for Java. In *OOPSLA 2000*, pages 264–280, 2000.
- [21] Frank Tip and Jens Palsberg. Scalable propagation-based call graph construction algorithms. In *OOPSLA 2000*, pages 281–293, 2000.
- [22] Raja Vallée-Rai, Etienne Gagnon, Laurie J. Hendren, Patrick Lam, Patrice Pominville, and Vijay Sundaresan. Optimizing Java bytecode using the Soot framework: Is it feasible? In *CC 2000*, pages 18–34, 2000.
- [23] Ian H. Witten and Eibe Frank. *Data Mining: Practical Machine Learning Tools and Techniques with Java implementations*. Morgan Kaufmann Publishers, 2000.