# The Obfuscation Executive[*]

Kelly Heffner and Christian Collberg

Department of Computer Science
The University of Arizona
{kheffner,collberg}@cs.arizona.edu

**Abstract.** Code obfuscations are semantics-preserving code transformations used to protect a program from reverse engineering. There is generally no expectation of complete, long-term, protection. Rather, there is a trade-off between the protection afforded by an obfuscation (i.e. the amount of resources an adversary has to expend to overcome the layer of confusion added by the transformation) and the resulting performance overhead.

In this paper we examine the problems that arise when constructing an *Obfuscation Executive.* This is the main loop in charge of a) selecting the part of the application to be obfuscated next, b) choosing the best transformation to apply to this part, c) evaluating how much confusion and overhead has been added to the application, and d) deciding when the obfuscation process should terminate.

## 1   Introduction

A *code obfuscator* is a tool which—much like a code optimizer—repeatedly applies semantics-preserving code transformations to a program. However, while an optimizer tries to make the program as fast or as small as possible, the obfuscator tries to make it as *incomprehensible* as possible. Obfuscation is typically applied to programs in order to protect them from being reverse engineered or to protect a secret stored in the program from being discovered. In this paper we will describe the design of an *Obfuscation Executive* (OE), implemented within the SANDMARK [3] software protection research tool. The OE is the overall loop that applies obfuscation algorithms to parts of the program to be protected. In many ways the OE functions similar to a compiler's optimization pass: it reads and analyzes an application and repeatedly applies semantics-preserving transformations until some termination condition has been reached. Ideally, the executive should be able to pick an optimal set of transformations and an optimal set of program parts to obfuscate. The only necessary user interaction should be to indicate to the tool what "optimal" means: i.e. how much execution overhead the user can accept, how much obfuscation he wants to add, and which parts of the application are security- or performance-critical. The two major issues with this process is the order in which transformations should be applied (the

---

"phase-ordering-problem") and how to decide that the process should terminate. In the case of a code optimizer the order of transformations is usually fixed. The optimizer typically terminates when there are no more changes to the application or when all transformations have been tried at least once. As we will see from this paper, in the case of the OE neither phase-ordering nor termination is this simple.

For the sake of brevity, we refer to [5] for an introduction to code obfuscation and its uses. Obfuscating transformations are characterized by their *potency* (the amount of confusion they add), their *resilience* (the extent to which they can be undone by a *de-obfuscator*), and their *cost* (the performance penalty they incur on the obfuscated application) [6].

## 2   The Obfuscation Loop

At the heart of any OE is a loop that chooses a part of the application to obfuscate, chooses an appropriate obfuscating transformation from a pool of candidate algorithms, and then applies the transformation. After the transformation has completed, the loop computes how much the code has changed and decides if the process should continue. Unfortunately, there are a number of complications. First of all, neither the amount of protection we would like to achieve nor the amount of overhead we can accept are uniform over the application. Some subroutines may be performance-critical, others not. Some subroutines may be security-critical, others not. Secondly, given any non-trivial set of obfuscating transformations there will be restrictions on the order in which these should be applied. The reason is that an obfuscating transformation *destroys* structures in the application. This makes the obfuscated application more difficult to analyze, and, as a result, it might not be possible to apply any further transformations. Finally, not all obfuscations can be applied to all application objects. For example, obfuscations which rename classes [9] cannot be applied to classes that will be loaded dynamically by name. Also, if we are using obfuscations to hide a particular structure in the application (such as a watermark or a set of cryptographic keys) we cannot apply transformations that will destroy these structures.

### 2.1   Transformation Dependencies

SandMark's OE understands six types of dependencies between obfuscations and tools such as software watermarkers that use them. These are pre/post-suggestions, pre/post-requirements, and pre/post-prohibitions:

**Pre-/postrequirement:** Assume a software watermarking algorithm (such as CT [4]) that embeds the watermark in a data structure in the application. To simplify implementation and debugging the watermarker creates a class `Watermark` that contains the code for building the mark: ⌜ `class Watermark { Watermark left, right; void createGraph() { ... } }` ⌝ A call to the method `createGraph()` is embedded into the application. Obviously,

this is not stealthy. Therefore the watermarking algorithm requires that an inlining transformation and a name obfuscating transformation be run after the watermarking algorithm. This is a *postrequirement* dependency.

**Pre-/postprohibition:** Assume an obfuscating transformation `MergeArrays` that performs alias analysis to detect the location of two arrays to merge. This algorithm should not be run after any algorithm that makes alias analysis difficult. Collberg [6] presents such algorithms. This type of dependency is called a *preprohibition.*

**Pre-/postsuggestion:** Suggestions are similar to requirements in the resulting language of transformations that they allow. However, while breaking a requirement will put the program in a corrupt state or make a software watermark obvious, a suggestion is just a hint to the OE by the obfuscation author that certain transformations work well together.

In our current SANDMARK implementation each obfuscation and watermarking algorithm specifies the effects that it may have on the code. It also specifies *properties* of other algorithms that are postrequired, presuggested, etc. For example, a method splitting transformation `MethodSplit` might list `OBFUSCATE_METHOD_NAMES` as a postsuggestion, indicating to the OE that some algorithm (any one will do) that has this property should be run after `MethodSplit`. In general, to fulfill a requisite dependency **one** algorithm with the specified property must be run; to fulfill a suggestion dependency any **one** algorithm with the specified property could be run; and to fulfill a prohibition dependency **no** algorithm with a specified property should be allowed to run.

Simple obfuscators do not need to specify many dependency relationships with other obfuscators. However, when using obfuscation with software watermarking, dependencies are necessary to make sure that the obfuscations successfully camouflage the watermark without destroying the structures that it is embedded in. From a software design standpoint, the dependency framework also allows for more complex obfuscators to be created modularly and without concern for the transformations performed around it.

## 3   Modeling Dependencies

Our goal is to construct an OE algorithm to honor transformation dependencies and to find the "optimal" set of transformations to apply to the "optimal" set of application objects. In fact, we are interested in constructing *families* of such algorithms, to be targeted at the many and varied applications of obfuscation discussed in [5]. The most important result in this paper is the design of a model which encodes transformation dependencies, obfuscation potency and performance overhead, as well as the desired level of obfuscation and overhead of each application object. The model is based on weighted finite state automata and can be used as the basis for many on-line and off-line algorithms. Let us first consider an example with three transformations $A$, $B$, and $C$. $A$ postprohibits a property that $C$ has, $A$ postrequires a property that only $B$ has, and $B$ prerequires a property that only $C$ has:The order in which these transformations

**A postrequires B**    **A prerequires B**    **A postprohibits B**    **A preprohibits B**

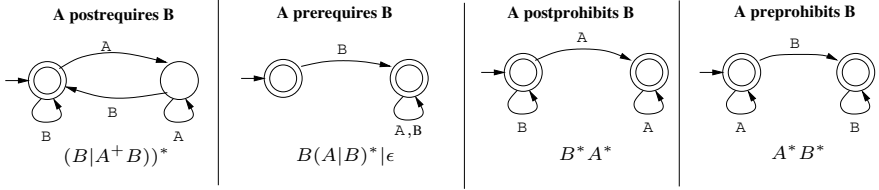$(B|A^+B))^*$        $B(A|B)^*|\epsilon$        $B^*A^*$        $A^*B^*$

**Fig. 1.** An FSA for each dependency type.

can be run is $CAB$, $CBAB$, $CBBAB$, or, more generally, any sequence that matches the regular expression $\epsilon|C(C|B)^*(A^+B)^*$.

This observation leads us to a model where the dependencies between transformations are represented by a finite state automaton (FSA). The language accepted by this machine contains all the possible sequences of transformations that can be executed. Given this model, the design of a new OE reduces to the problem of constructing an algorithm that chooses a finite subset of the (typically) infinite set of strings generated by the FSA. The heuristics of such algorithms will make use of FSA *edge weights* representing the "goodness" of traversing each edge.

We will next describe how the FSA is built, then how edge weights are computed, and, finally, in Section 4 describe an on-line OE algorithm that makes use of the FSA model.

### 3.1   Building the FSA

Each type of dependency has an equivalent regular language. It suffices to consider prohibition and requirement dependencies since suggested dependencies can be modeled by modifying the FSA edge weights. This will be shown in Section 3.2.

Figure 1 shows the four FSAs that correspond to each dependency type. The figures show transitions only for the two transformations involved in the dependency; for any transformation that is not in the dependency, the transition is just a self-loop.

To build the regular language for the entire set of obfuscating transformations we take the intersection of the languages from each dependency and $\Sigma^*$. The resulting language is the set of all possible sequences in which the transformations could be applied. To model the fact that dependencies apply to properties of transformations, rather than transformations themselves, we simply replace a single transformation in Figure 1 with all of the transformations that have a particular property.

Consider the running example in Figure 2 which describes five obfuscating transformations $A$, $B$, $C$, $D$, and $E$ with three properties $p_1$, $p_2$, and $p_3$. Figure 3(a) shows the models for the four dependencies in the example. Taking the intersection of the languages represented by these FSAs we get the FSA

```
class c1{
  m1(){...}
  m2(){...}
}
class c2{
  m(){...}
}
```

| Transfor-mation | Obfuscation | | | Prop-erties | Required | | Prohibited | |
|---|---|---|---|---|---|---|---|---|
| | Level | Potency | Overhead | | Pre | Post | Pre | Post |
| A | Method | 1.0 | 0.9 | $p_1$ | $p_2$ | | | |
| B | Method | 0.5 | 0.3 | $p_2, p_3$ | | $p_1$ | | |
| C | Class | 0.1 | 0.4 | $p_1, p_3$ | | | | |
| D | Class | 0.2 | 0.2 | $p_2$ | | | $p_3$ | |
| E | App. | 0.01 | 0.1 | $p_3$ | | | | $p_1$ |

**Fig. 2.** A running example. $A$ and $B$ are method level obfuscators, $C$ and $D$ are class level obfuscators, and $E$ applies to an entire application.
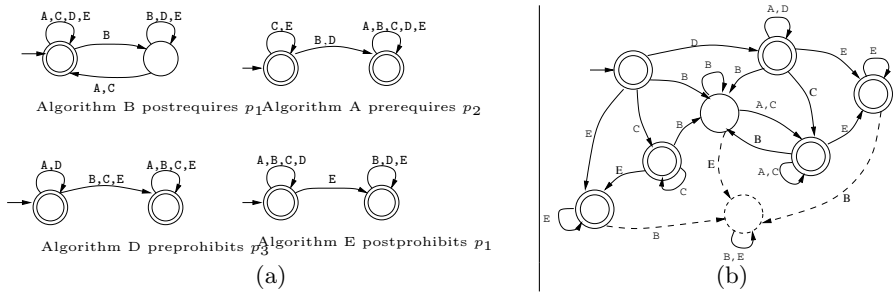


**Fig. 3.** The transformation properties from the running example in Figure 2 produces the four FSAs in (a). Taking the intersection of the generated languages produces the FSA in (b). Note that no paths exist to an accepting state from some states (dashed in (b)) so those nodes would be removed.

in Figure 3(b), which models every possible candidate sequence allowed by the dependencies.

To make full use of the finite state machine model, we must integrate the idea that transformations are run on application objects, not always the entire application. By running a transformation on one application object, possibly fulfilling prerequirements and prohibiting other transformations, the result affects only a subset of objects. In order to keep track of these object-level changes, the target of the transformation must be included with each transformation in the sequence. Thus, each symbol in our alphabet for the sequence becomes an ordered pair (*transformation*, *target*).

We should also note that modifications to single application objects do not just affect that object. Obfuscating a method may affect not only the class that the method is in, but the entire application. For simplicity we will assume that when a transformation is run on an object $x$, the effect is spread to all objects that contain $x$, and all objects that $x$ contains. We will call this the *range* of an object. This is a conservative approximation on the real spread. A less coarse approximation would be advantageous but more difficult to compute.

During FSA construction the states represent sets of application objects. We will refer to the set of objects for a given state $q$ as $s(q)$. We will refer to the set of objects that is in the range of an application object $x$ as $r(x)$.

First we will construct the FSA for a prerequirement dependency, where a transformation $T$ prerequires a property $p$. We define an FSA $(Q, \Sigma, \delta, q_0, F)$ with the following properties: $Q$, the set of states, is composed of the power set of the set of all application objects that are a target of $T$; $\Sigma$ is the set of ordered pairs $(T, x)$ where $T$ is a transformation and $x$ is a target application object for that transformation; The transition function $\delta$ is given in Figure 5; $q_0$ is the state such that $s(q_0) = \emptyset$; $F = Q$ is the set of accepting states. Figure 6 shows the partial FSA for the preprohibition of property $p_3$ before transformation $D$.

Next, we construct the FSA for a postrequirement dependency, where transformation $T$ postrequires property $p$. The FSA is identical to the previous one, except for the transition function in Figure 5 and that $F = \{q_0\}$. The transition functions for FSAs for preprohibition and postprohibition dependencies are similar and shown in Figure 5.

## 3.2   Building the Probabilistic FSA

The FSA generates a language of strings of (*obfuscation*, *target*) tuples. These strings represent a series of obfuscations to run on the application. Obviously, some strings are more desirable than others in that they result in more highly obfuscated programs with lower performance penalty. To capture this, we create a probabilistic FSA by giving each edge an *edge weight*. In an edge $a \xrightarrow{(T,x),w} b$, $w$ represents the "goodness" of applying transformation $T$ to application object $x$ when the OE is in state $a$. In Section 4 we will show how this model allows for a very simple, yet effective, on-line OE algorithm.

The tuple

$$\langle Potency(T), Degradation(T), ObfLevel(x), PerfImport(x) \rangle$$

forms the FSA edge weight $weight(T, x)$, where each element is a real number in the range $[0, 1]$. $Potency(T)$ measures the obfuscation potency of $T$. It is computed by running each obfuscation on a set of benchmarks and computing the change in software complexity. See Section 3.3. $Degradation(T)$ measures the performance degradation of $T$. This is computed by running each obfuscation on a set of benchmarks and computing the change in execution time. See Section 3.3. $ObfLevel(x)$ is the desired obfuscation level of application object $x$, as assigned by the user. $PerfImport(x)$ is the importance of performance of application object $x$, a combination of user assignment and profiling data. Our model does not specify how a particular OE will make use of the weight tuple. In our current implementation the tuple is mapped down to a single real number which represents the overall goodness of choosing a particular edge. This is explained in Section 3.3. We will next show how to estimate $Potency(T)$ and $Degradation(T)$.

## 3.3   Computing Edge Weights

Each obfuscating transformation $T$ is assigned a real number (in the range $[0, 1]$) $Potency(T)$ that represents the relative potency of the transformation in comparison with the rest of the transformations known to the obfuscator. $Potency(T)$

(a) $\Delta(P, M, T) = \frac{\sum_{m \in M} |m(P) - m(T(P))|}{|M|}$ (b) $Potency(T) = \frac{\Delta(M,T) - \Delta_{\min}(M)}{\Delta_{\max}(M) - \Delta_{\min}(M)}$

(c) $\Gamma(P, T) = \frac{\max(\text{time}(T(P)) - \text{time}(P), 0)}{\text{time}(P)}$ (d) $Degradation(T) = \frac{\Gamma(T) - \Gamma_{\min}}{\Gamma_{\max} - \Gamma_{\min}}$

(e) $Fold(w) = Potency(T) \cdot ObfLevel(x) \cdot (1 - PerfImport(x)) \cdot (1 - Degradation(T))$

**Fig. 4.** Formulas for computing edge weights.

is determined by calculating a set of software engineering metrics, $M$, on sample programs before and after obfuscation and taking the average of the change in those metrics (see Section 4.1 for discussion on what metrics were used).

The change in a set of metrics $M$ for a program $P$ on a transformation $T$ is given by Figure 4(a) where $m(P)$ is a software metric calculated on $P$ and $T(P)$ is $P$ obfuscated by transformation $T$. We get $\Delta(M, T)$ by averaging $\Delta(P, M, T)$ over all benchmark programs.

To compute the obfuscation potency for a transformation $T$ we normalize using the highest and lowest changes in the metrics ($\Delta_{\max}(M)$ and $\Delta_{\min}(M)$) over all of the obfuscations, yielding the formula in Figure 4(b). Since we considered obfuscation as simply a change in metrics, rather than raising or lowering the metrics, we chose to calculate $\Delta(P, M, T)$ as given by Figure 4(a). Another method of selecting metrics, such that more obfuscation mapped to higher metric values would eliminate the absolute value bars from the calculation.

Each obfuscating transformation $T$ is assigned a real number $Degradation(T)$ (in the range $[0, 1]$) that represents $T$'s expected performance hit. $Degradation(T)$ is calculated by running every transformation $T$ on a set of "priming" applications. These could either be a set of benchmarks such as the SpecJVM, or the application to be obfuscated itself. $\Gamma(P, T)$ is the raw performance degradation of $T$ on program $P$, shown in Figure 4(c). We get $\Gamma(T)$ by averaging over all benchmark programs. To compute $Degradation(T)$ we normalize using the highest and lowest performance hits ($\Gamma_{\max}$ and $\Gamma_{\min}$), yielding the formula in Figure 4(d).

Most simple OE algorithms will want to fold the weight tuple $w = weight(x, T)$ into a single real number $Fold(w)$ to represent the goodness of choosing a particular edge. See Figure 4(e). The probability of taking an edge $a \xrightarrow{(T,x), Fold(w)} b$ is thus proportional to how potent the transformation $T$ is and how important it is to obfuscate application object $x$. It is *inversely* proportional to how performance critical $x$ is and how much $T$ is expected decrease its performance.

## 4   Algorithms

Once the probabilistic FSA has been constructed, the model is used to find an effective, yet not optimal, obfuscation sequence. Our broad definition for an optimal obfuscation sequence is a sequence such that the application has maximal obfuscation, has minimal performance degradation, and is a minimal

| | requirement | prohibition |
|---|---|---|
| **pre** | $q'$, if $q \neq q'$, $t$ has the property $p$, and $s(q) + r(x) = s(q')$, or<br>$q$, if $t = T$ and $x \in s(q)$, or<br>$q$, if $t$ has property $p$ and $r(x) \subseteq s(q)$, or<br>$q$, if $t \neq T$ and $t$ does not have property $p$ | $q'$ if $t$ has the property $p$, and $s(q) + r(x) = s(q')$, or<br>$q$ if $t = T$ and $x \notin s(q)$, or<br>$q$ if $t$ has property $p$ and $r(x) \subseteq s(q)$, or<br>$q$ if $t \neq T$ and $t$ does not have property $p$ |
| **post** | $q'$ if $t = T$ and $s(q) + r(x) = s(q')$, or<br>$q'$ if $t$ has the property $p$, and $s(q) - r(x) = s(q')$, or<br>$q$ if $t = T$ and $r(x) \subseteq s(q)$, or<br>$q$ if $t$ has property $p$ and $r(x) \cap s(q) = \emptyset$, or<br>$q$ if $t \neq T$ and $t$ does not have property $p$ | $q'$ if $t = T$ and $s(q) + r(x) = s(q')$<br>$q$ if $t = T$ and $T$ does not have property $p$ and $x \in s(q)$, or<br>$q$ if $t$ has property $p$ and $r(x) \cap s(q) = \emptyset$, or<br>$q$ if $t \neq T$ and $t$ does not have property $p$ |

**Fig. 5.** Transition functions $\delta(q, (t,x))$.

sequence. Here, we describe a simple on-line algorithm to determine a sequence. Given the probabilistic FSA and the *Fold* formula from Section 3.3, the on-line algorithm computes the obfuscation sequence and performs the associated obfuscating transformations.

The algorithm performs a random walk of the nodes of the FSA, starting in the start node. Each iteration selects an outgoing edge $e$ from the current node $S$, performs the corresponding obfuscating transformation, and updates the edge weights to reflect the changing obfuscation level. The probability of a particular outgoing edge being chosen is proportional to its weight. As the desired obfuscation level of each application object approaches zero, the weights of the edges that represent obfuscating that application object also approach zero. Edges with a zero weight are removed from the FSA. The loop terminates when there are no available edges out of the current, accepting, state.

This algorithm is a random walk of the FSA, using the edge weights to guide traversal. While this algorithm will not yield an optimal obfuscation sequence, it will produce a sequence that obfuscates heavily with acceptable performance degradation. Furthermore, since this algorithm produces a very random obfuscation sequence, attacks against the obfuscated code are difficult. Given a list of the obfuscations available to the loop, the attacker still does not know the subsets of obfuscation that the application objects have been obfuscated with, nor the order in which the obfuscations have been applied. We can use this loop to obfuscate a fingerprinted application each time it is sold. Each copy will become entirely different, allowing us to protect against collusive attacks.
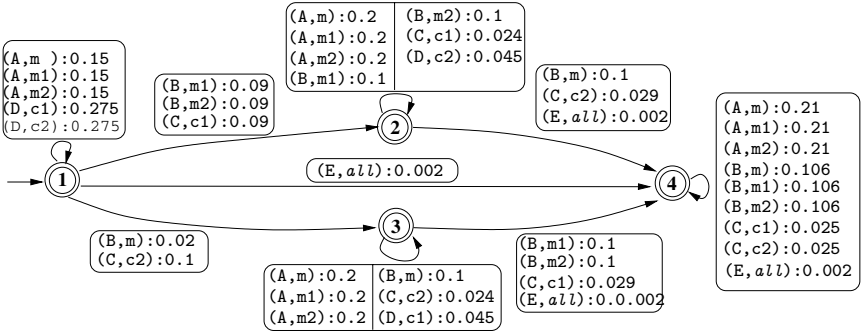
**Fig. 6.** The FSA that models the preprohibition of property $p_3$ before transformation $D$. Each edge is labeled with a set of tuples $(T, x) : w$ where $T$ is the obfuscation transformation to run on application object $x$ and $w$ is the weight of tuple (see Section 3.2). For each node the weights of the outgoing edges sum to 1, which allows the FSA model to be used as a probabilistic FSA.

Consider again the running example from Figure 2 and the FSA model shown in Figure 6. The weights are shown as they would be computed before the first iteration of the loop. We assume that the obfuscation level of all objects is 1 and that performance importance is 0. On the first iteration we randomly choose to move from state 1 (the start state) to state 3 by running algorithm C on class c2. This will cause c2 to be obfuscated, lowering its remaining *ObfLevel*(c2). During the next iteration any edge with c2 as its target will have a lower weight, lowering its probability to be obfuscated again. Note that moving to state 3 eliminates the possibility of ever running D on c2. This is because algorithm D preprohibits any algorithm with property $p_3$, such as C.

### 4.1   Implementation

The probabilistic FSA constructed in Section 3 provides a clean model for the dependencies between transformations. It has also allowed us to construct the simple random walk OE algorithm above. However, a straight-forward implementation of these ideas turns out to be impractical. The reason is that for sets of transformations with very few dependencies the size of the FSA grows exponentially. The solution to this problem is to lazily build and walk the FSA concurrently.

In our implementation, the metric set $M$ used for *Potency*(T) is computed based on a suite of standard software complexity metrics [8], including the standard ones proposed by McCabe, Halstead, Chidamber, Harrison, Munson, and Henry.

To compute *Degradation*(T) we use a large suite of standard Java benchmarks including SPEC JVM98 (www.specbench.org/osg/jvm98) and the Ashes test suite (www.sable.mcgill.ca/ashes).
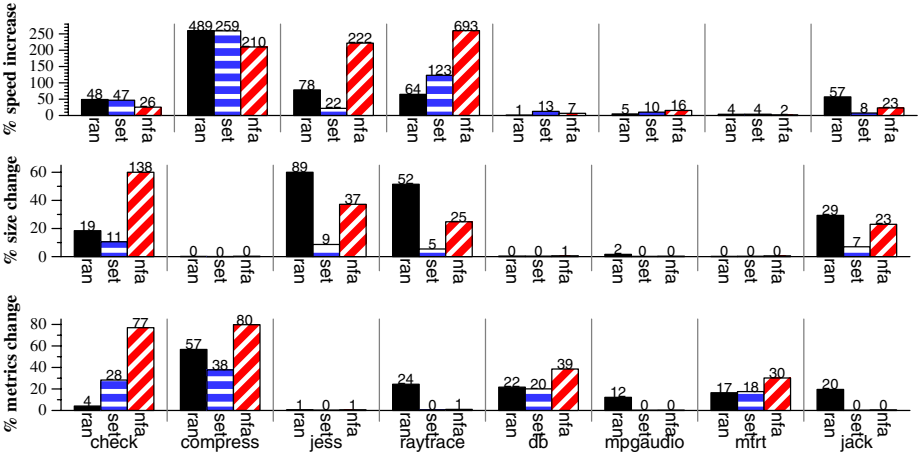
**Fig. 7.** Results from implemented obfuscation loops on SpecJVM.

Our current implementation supports three OEs: a simplistic set-based model (described in Section 4.2 below), the complete FSA-based model, and the lazy FSA model. All are on-line algorithms. We are currently exploring an off-line OE which uses the probabilistic FSA to compute an "optimal" obfuscation sequence ahead of time.

## 4.2   Evaluation

To evaluate the FSA-based OE algorithm we compare it to two simpler OE algorithms. Each algorithm was run on SpecJVM benchmarks with the desired obfuscation level for each object maximized and no user input about application hotspots. Figure 7 shows the result of running the three OEs. For each benchmark we show the change in code size, execution time, and software complexity metrics.

**Random Select OE.** The first OE algorithm, Random, does not split the program into obfuscation objects but instead applies obfuscating transformations to the entire program with each pass. Random applies a few heuristics to manage transformation dependencies, specifically enforcing the prerequisite, preprohibition, and postprohibition rules, but makes no effort to ensure postrequisites can be fulfilled. After determining which transformations cannot be applied to the program, a transformation is chosen at random from the remaining candidates. If the executive detects that the program has been transformed into a corrupt state, it simply halts. After each obfuscation, it computes the change in complexity metrics and halts when the desired amount of obfuscation has been applied.

The Random implementation fails to find an appropriate cost-benefit trade-off for obfuscation, notably in check, compress, and jess where the change in

metrics is low in comparison to the FSA for a large trade-off in size and speed. It is important to note that the random implementation applies each obfuscation to the entire program. The results of the Random OE shows that simply running all of the obfuscations on the program randomly does not produce an acceptable amount of obfuscation; there is a need for an intelligent OE.

**Set-Based OE.** The second OE algorithm implements a set-based model. In this algorithm, each application object is associated with a set of obfuscation candidates, transformations that are allowed to be run on that object. In each iteration the loop uses a set of heuristics to trim the candidate set to remove algorithms that have been disqualified by running the previous obfuscation. The target application object is chosen based on the desired obfuscation level for each object and the amount of obfuscation already applied to the object. Another set of heuristics is used to extract a subset of obfuscations that can be run and in most cases will not lead to a corrupted state where required dependencies cannot be fulfilled. A candidate is then chosen at random from this subset. The set heuristics fail to detect cases such as two obfuscating transformations conflicting with each other, yielding only the empty-string as a valid obfuscation sequence. The set-based model halts when each application object reaches the desired obfuscation level.

The set-based model fails to impact the complexity metrics over the entire program, while still incurring the cost of speed and size for the attempted obfuscations. Presumably, this is due to the fact that the set-based implementation chooses the next obfuscation target in isolation to the entire application, without considering which (*transformation, object*) pair will yield the most obfuscation overall.

**FSA-Based OE.** As expected, the FSA produces a large increase in software complexity, at the cost of program size and performance. The lack of user input about hotspots is most evident in the cases of extreme performance hit, like `raytrace`. It would also be useful to re-calculate profiling data between each iteration of the obfuscation loop, with the obvious performance implications.

Overall, the numbers show that there is much to be gained over random OEs. We see that the FSA-based model had the largest average change to engineering metrics at 28.5% (compared to 19.6% and 13% for random and set) however at the cost of a larger speed penalty (50% speed increase compared to an increase for random and set of 93% and 60%). The important thing to note is that the other two models show a large downgrade to the software (random has an average size increase of 24%, very close to the FSA-based change of 27%) without affecting the metrics by the same degree as the FSA-based OE (the set based model had an average speed increase of 60% with only an average change in metrics of 13%). More analysis of individual obfuscations and their effect on speed, size, and complexity, using complexity measures to drive the direction of obfuscation, and integrating knowledge of the program designer will lead to a more intelligent OE.

The FSA-based model is superior to the other approaches. It elegantly handles the obfuscation dependencies, without the use of conservative restrictions to guarantee that the loop will not go into a corrupted state. The FSA-based model can yield several algorithms using different analyses, whereas the set-based approach is tied to the trimming algorithm and can only be adjusted in the way that the obfuscation target and candidate are chosen. In addition, modeling obfuscation sequences as weighted members of a regular language provides a straightforward solution for deriving a *family* of effective transformation sequences for use in artificial diversity.

## 5    Discussion and Summary

To the best of our knowledge, Collberg et al. [5] is the first description of an OE. This algorithm does not take into account restrictions on transformation ordering. Wroblewski [12] describes an OE for x86 machine code. This OE applies obfuscations in a single pass over the program using a hardcoded sequence of transformations. Lacey et al.  [7] presents an algorithm which decides whether applying one optimizing transformation to a piece of code will prevent another one from being applied.

There are few theoretical results related to obfuscation. Barak et al.  [2] shows that there exist programs that cannot be obfuscated. Appel [1] shows deobfuscation to be NP-easy. The proof idea is based on a simple algorithm which nondeterministically guesses the original program $S$ and obfuscation key $K$. The obfuscating transformation is then run over $S$ and $K$, verifying that the result is the obfuscated program. To use this algorithm to defeat our random walk obfuscation loop, $K$ must include the seed to our random number generator. However, Appel's algorithm is only valid for obfuscating transformations that are injective; transformations such as name obfuscation or instruction reordering cannot be reversed using Appel's method.

Determining the order in which to run obfuscation and watermarking algorithms is similar to the phase ordering problem for optimizations in compilers. In the case of optimizers, the phases consist of code analyses and optimizing transformations. While most compilers hardwire the order of the optimizations, some work [11,10] has been done on automated phase selection and ordering.

The widely differing uses of software obfuscation have led us to the very general probabilistic FSA model of Section 3. The model encodes all valid obfuscation sequences as well as the goodness of any particular sequence. Two on-line OE algorithms have been implemented that makes use of this model. These algorithms compare favorably to an algorithm based on a simplistic set-based OE model. The SANDMARK framework can be downloaded from `http://sandmark.cs.arizona.edu`. SANDMARK consists of 130,000 lines of Java, and includes 39 obfuscation, 17 watermarking, 3 OE algorithms, and several tools for automatically and manually analyzing and attacking software protection algorithms.

# References

1. A. Appel. Deobfuscation is in np.
   `www.cs.princeton.edu/~appel/papers/deobfus.pdf`.
2. B. Barak, O. Goldreich, R. Impagliazzo, S. Rudich, A. Sahai, S. Vadhan, and
   K. Yang. On the (im)possibility of software obfuscation. In *Crypto01*, pages 1–18,
   2001. LNCS 2139.
3. C. Collberg, G. Myles, and A. Huntwork. Sandmark - a tool for software protection
   research. *IEEE Security and Privacy*, 1(4):40–49, 2003.
4. C. Collberg and C. Thomborson. Software watermarking: Models and dynamic
   embeddings. In *POPL'99*, San Antonio, TX, Jan. 1999.
5. C. Collberg, C. Thomborson, and D. Low. A taxonomy of obfuscating transfor-
   mations. Technical Report 148, Department of Computer Science, University of
   Auckland, July 1997.
6. C. Collberg, C. Thomborson, and D. Low. Manufacturing cheap, resilient, and
   stealthy opaque constructs. In *POPL'98*, San Diego, CA, Jan. 1998.
7. D. Lacey and O. de Moor. Detecting disabling interference between program trans-
   formations. `citeseer.nj.nec.com/464977.html`.
8. S. Purao and V. Vaishnavi. Product metrics for object-oriented systems. *ACM
   Comput. Surv.*, 35(2):191–221, 2003.
9. H. P. V. Vliet. Crema — The Java obfuscator.
   `web.inter.nl.net/users/H.P.van.Vliet/crema.html`, Jan. 1996.
10. D. Whitfield and M. L. Soffa. An approach to ordering optimizing transformations.
    In *PPOPP'90*, pages 137–146, 1990.
11. D. Whitfield and M. L. Soffa. Automatic generation of global optimizers. In
    *PLDI'91*, pages 120–129, 1991.
12. G. Wroblewski. *A General Method of Program Code Obfuscation*. PhD thesis,
    Wroclaw University, 2002.