

ML for ML: Learning Cost Semantics by Experiment

Ankush Das^(✉) and Jan Hoffmann

Carnegie Mellon University, Pittsburgh, USA
ankushd@cs.cmu.edu

Abstract. It is an open problem in static resource bound analysis to connect high-level resource bounds with the actual execution time and memory usage of compiled machine code. This paper proposes to use machine learning to derive a cost model for a high-level source language that approximates the execution cost of compiled programs on a specific hardware platform. The proposed technique starts by fixing a cost semantics for the source language in which certain constants are unknown. To learn the constants for a specific hardware, a machine learning algorithm measures the resource cost of a set of training programs and compares the cost with the prediction of the cost semantics. The quality of the learned cost model is evaluated by comparing the model with the measured cost on a set of independent control programs. The technique has been implemented for a subset of OCaml using Inria’s OCaml compiler on an Intel x86-64 and ARM 64-bit v8-A platform. The considered resources in the implementation are heap allocations and execution time. The training programs are deliberately simple, handwritten micro benchmarks and the control programs are retrieved from the standard library, an OCaml online tutorial, and local OCaml projects. Different machine learning techniques are applied, including (weighted) linear regression and (weighted) robust regression. To model the execution time of programs with garbage collection (GC), the system combines models for memory allocations and executions without GC, which are derived first. Experiments indicate that the derived cost semantics for the number of heap allocations on both hardware platforms is accurate. The error of the cost semantics on the control programs for the x86-64 architecture for execution time with and without GC is about 19.80% and 13.04%, respectively. The derived cost semantics are combined with RAML, a state-of-the-art system for automatically deriving resource bounds for OCaml programs. Using these semantics, RAML is for the first time able to make predictions about the actual worst-case execution time.

1 Introduction

Motivated by longstanding problems such as performance bugs [32], side-channel attacks [10, 31], and to provide development-time feedback to programmers, the programming language community is developing tools that help programmers understand the resource usage of code at compile time. There has been great progress on automatically determining loop bounds in sequential

C-like programs [13, 15, 22, 35], deriving, solving recurrence relations [4, 6, 9, 19], and automating amortized analysis [23–25]. There exist several tools that can automatically derive loop and recursion bounds, including SPEED [21, 22], KoAT [13], PUBS [5], Rank [7], ABC [11], LOOPUS [35, 38], C4B [14], and RAML [23, 24].

Most of these resource analysis tools use high-level cost models, like number of loop iterations and function calls, and it is often unclear how the derived bounds relate to the machine code executing on a specific hardware. To make the connection, one has to take into account compilation, hardware specific cache and memory effects, instruction pipelines, and garbage collection cycles. While there exist tools and techniques for analyzing low-level assembly code to produce worst-case execution time bounds for concrete hardware [37], they are limited in their expressive power, as analyzing assembly code is a complicated problem.

In this article, we propose a novel technique to derive cost models that can link high-level resource bounds to the execution of low-level code. We present a simple operational cost semantics for a subset of OCaml [20] that have been learned using standard machine learning techniques like linear regression. The resources we are considering are heap allocations, execution time without garbage collection (GC), and execution time with GC. The subset of OCaml we are considering is purely functional and includes lists, tuples and pattern matching. However, the technique is also applicable to programs with side effects.

To learn a cost semantics, we first define an operational big-step semantics that assign a parametric cost expression to a well-formed expression. This cost expression is parametric in (yet unknown) constants that correspond to high-level constructs in OCaml. The assumption is that the number of executions of each of these constructs constitutes the majority of the execution time of the expression. We keep track of the number of executions of each of these constructs in the cost semantics, which has been implemented in an OCaml interpreter. Our semantics then models the execution time of the program as a linear sum of the number of executions of each construct. The (unknown) coefficients of this linear function intuitively represent the execution time of each construct.

We then determine the average values of the coefficients on a specific hardware by experiment. We carefully select a set of relatively simple training programs and measure the median execution time of these programs on the hardware of interest. To this end, each program is executed with the OCaml native code compiler 500 times on an Intel x86-64 and a ARM 64-bit v8-A platform. We then apply machine learning techniques, such as linear regression [30], on the linear functions obtained by running the cost semantics to determine the constant costs of the constructs by fitting the prediction of the cost semantics to the measured costs. We measure the execution time using the Unix library in OCaml, which is hardware independent. We measure the number of allocation by relying on the OCaml GC library, which is again hardware independent. As a result, our approach is completely hardware independent and can be easily extended to different architectures, as we demonstrate in this work.

Of course, the execution time of, say, an addition cannot be described by a constant. In fact, it can vary by a large margin depending on whether the arguments are stored on the stack or in a register. Similarly, a cons operation

can be costly if one of the arguments is not in the cache has to be fetched from memory. So the constants that we learn in our experiment will represent roughly the *average cost* of the operations on the training programs.

Once we have learned these cost coefficients for a specific hardware and resource metric, we validate our cost with control (or test) programs, retrieved from the standard library, an OCaml online tutorial, and local OCaml projects. Each control program is first compiled and executed on the hardware and the median execution cost is measured in the same way we did for training programs. The program is then run on the interpreter to obtain the parametric linear cost function. By plugging in the learned coefficients, we get a prediction for the execution time or memory usage. We compare the predictions of the cost semantics with the median cost, and report the percentage error on test programs. We use the median instead of the mean because it is more resilient against outliers which are often caused by context switches in the OS.

The result of the experiments with the control programs are surprisingly encouraging. We precisely learn the amount of memory that is allocated by each construct. For execution time of programs that do not trigger GC, the error of our model is up to 43%, for all but one program.

In memory intensive programs, the impact of garbage collection cycles on the execution time is significant. So, we adapt our cost semantics to account for the time taken by the GC. We make two simplifying assumptions to model the GC time. One of them is that the time taken by each GC cycle is a constant and the other is that each GC cycle starts with a full heap, and ends with an empty heap. These assumptions, as we will see in the experiments and the OCaml documentation, are quite close to the collections of the *minor heap*. To model this behavior, we combine the cost semantics for memory allocations and the cost semantics for programs without GC. Since the GC cycle occurs periodically when the minor heap is full, we can predict the number of minor GC cycles in the lifetime of a program using the allocation semantics. To determine the time needed for a minor garbage collection, we just measure the median GC time taken by a GC cycle for the training programs.

The main application of our cost semantics is the integration into Resource Aware ML (RAML), a state-of-the-art tool for automatic resource analysis. Using the semantics for execution time on x86, RAML can automatically derive worst-case bounds for many functions that are close to the measured execution time of the compiled code. Our results are precise enough, to statically determine the faster versions of different implementations of list append, Sieve of Eratosthenes, and factorial.

2 Method and Experimental Setup

In this section, we describe our experimental setup and training method. The main hypothesis, going into this experiment, is that the resource consumption of a program, whether time or memory, is a linear combination of the number of executions of each construct in the program. Moreover, the time (or memory) consumed by each construct is averaged out to be a constant. Hence, the execution time of a program is

$$T = \sum_{c \in \mathcal{C}} n_c T_c \quad (1)$$

where \mathcal{C} represents the set of constructs and n_c is the count of each construct during program execution, and T_c is the execution time of the respective construct. Clearly, these assumptions do not, in general, hold for most of the relevant platforms. Hence, we will analyze the error incurred by these simplifying assumptions on the execution time (or memory allocation) of a program.

Consider the following OCaml program, which computes the factorial.

```
let rec fact n = if (n = 0) then 1 else n * fact (n-1);;
(fact 10);;
```

In the above program, if we count the number of high level constructs, we get 10 function calls, 11 equality checks, 10 subtractions and multiplications and 1 “let rec” that defines the function. In our model the execution time of a program is the sum of the execution time of each construct multiplied by the number of times that construct is executed. For the above program, the total execution time is $11 * T_{FunApp} + 11 * T_{IntEq} + 10 * T_{IntSub} + 10 * T_{IntMult} + 1 * T_{letrec}$.

We are interested in the resources costs T_i that best approximate the actual cost. With this representative example, we describe our experimental setup.

Language Description. We have chosen a subset of OCaml as our modeling language. In this subset, we include the following program constructs: recursive functions, conditionals, boolean, integer and float comparisons and arithmetic, pattern matching and tuples. With this fairly general subset, we can write a variety of programs including list manipulation, matrix operations and other numeric programs, as we will demonstrate in our results. We chose OCaml as the source language for several reasons. For one, OCaml is a widely used language for functional programming which is quite efficient in practice. Moreover, we wanted to demonstrate that it is possible to define a practical cost semantics for a high-level functional language with a sophisticated compiler and automatic memory management. We assume that defining such a semantics would be easier for imperative programs, which are closer to assembly code.

A major obstacle when analyzing high-level languages is compiler optimization. The resource usage of the generated target assembly code depends greatly on the choices that are made by the compiler and cannot directly be derived from the original OCaml program. Hence, the cost semantics need to account for compiler optimizations. In our experience, we found two compiler optimizations with a significant impact on the execution time.

- *Tail Call Optimization* [36] - If the final action of a function body is a function call, it is optimized to a jump instruction. This is relevant for the cost semantics because a jump is faster than a call, hence we need two different costs for usual function calls and tail calls. Moreover, separating these costs in the semantics will later help us validate the claim that tail call optimization indeed reduces the execution time.
- *Function Inlining* [16] - OCaml compiler inlines functions as an optimization. Instead of accounting for inlining in our interpreter, we forced the compiler to

not inline any function when generating the native code. We will demonstrate the effect of this optimization when describing training programs. Conceptually, inlining is not a problem for our approach since it is possible to track at compile time which function applications have been inlined.

Training Algorithm. We formally describe the algorithm we use to learn the values of the constructs. Consider again our cost expression $T = \sum_{c \in \mathcal{C}} n_c T_c$. Essentially, we execute the native code obtained from the compiler to obtain the value T , and we execute the program on the interpreter to obtain the values n_c . Let there be P training programs and suppose we generate M instances of training data from each training program using the above method. We denote the count of each construct and execution time of the i -th instance of training data generated by j -th training program by $(n_c^{(i,j)})_{c \in \mathcal{C}}$ and $T_{(i,j)}$ respectively. Since we need to learn a linear model on T_c 's, linear regression is the natural choice of machine learning algorithm. A simple linear regression [33] would produce the following objective function.

$$S = \sum_{j=1}^P \sum_{i=1}^M \left(T_{(i,j)} - \sum_{c \in \mathcal{C}} n_c^{(i,j)} T_c \right)^2$$

where T_c are the unknowns that need to be learned. However, this approach is useful only when the error is additive, i.e. the error is independent of n_c . Unfortunately, in our case, each instruction has an inherent noise, which depends on the instruction, the operating system and synchronization overhead, measurement error and possibly other factors. So, as the number of instructions executed increases, the error in execution time also increases. Such an error is called *multiplicative*, and a solution by simple linear regression is skewed towards the constructs which have a higher cost, leading to inaccurate results. To overcome this problem, we need to normalize our objective function. We normalize the objective function by the sum of the execution time for each training program over all inputs. Hence, the new objective function is

$$S = \sum_{j=1}^P \sum_{i=1}^M \left(\frac{T_{(i,j)}}{S_j} - \sum_{c \in \mathcal{C}} \frac{n_c^{(i,j)}}{S_j} T_c \right)^2$$

where $S_j = \sum_{i=1}^M T_{(i,j)}$, i.e. the total execution time of the j -th training program. We learn the cost of each construct using the weighted linear regression technique. In addition to the above method, we also employ two other regression techniques. One is robust regression [34], where the objective function is the L_1 -norm, instead of the L_2 -norm (written as S_{RR} below).

$$S_{RR} = \sum_{j=1}^P \sum_{i=1}^M \left| \frac{T_{(i,j)}}{S_j} - \sum_{c \in \mathcal{C}} \frac{n_c^{(i,j)}}{S_j} T_c \right|$$

And the other is the non-negative least squares method [28], where the sum of squares is minimized under the constraint that all constants need to be

non-negative. We evaluate each algorithm, and compare the results obtained for each regression technique in Sect. 7.

Training Programs. The main goal of training is to learn appropriate values T_c for each program construct c in the language described above. Since we have $|\mathcal{C}|$ variables that we need to learn, all we need is at least $|\mathcal{C}|$ training programs to get a feasible cost semantics. However, there is a pitfall here that we want to avoid. Most of the typical machine learning algorithms suffer from the problem of overfitting, i.e. when the model learned is biased towards the training data, and performs poorly on the testing data. Specifically, in our case, the function call construct exists in all the training programs, hence, the learning algorithm overfits the data w.r.t. the cost for function call. Moreover, the regression algorithm is unaware of the fact that these costs need to all be positive. To overcome these issues, we need to linearly separate out the cost of each construct. To this end, we create one training program for each construct. Such a program has a significant count of one construct while being moderate in other constructs. For example, the training program for function call is

```
let id n = n;;
let rec fapp x = if (x = 0) then 0 else fapp (id (id (id (id (x-1)))));
```

If we don't account for function inlining, the function `id` gets inlined, and the above is treated as 1 application instead of 5. This massively impacts our training, and we obtain an incorrect cost for function application. Similarly, the training program for integer addition is

```
let rec fintadd x = if (x = 0) then 0 else x + x + x + x + fintadd (x-1);
```

Once we decided the training programs, we ran each training program with 20 inputs, ranging from 1000 to 20000. In this experiment, we have a total of 36 programs, and with 20 inputs for each program, we have a total of 720 training points for our linear regression algorithm. With this small training set, the regression techniques learn the cost model in less than 1 s. This training set might appear overly simplistic but our results show that this simple setup produces already surprisingly satisfying results.

Hardware Platforms. All of the above experiments have been performed on two separate platforms. One is an Intel NUC5i5RYH which has a 1.6 GHz 5th generation Intel Core i5-5250U processor based on the x86-64 instruction set. Another is a Raspberry Pi 3 which has a 1.2 GHz 64-bit quad-core ARM Cortex-A53 processor based on the ARM v8-A instruction set. We will report the results for both these platforms.

3 Operational Cost Semantics

In the following, we define the big-step operational cost semantics. A value environment V maps variables to values. An evaluation judgment $V \vdash e \Downarrow v \mid t$ denotes that in the environment V , the expression e evaluates to the value v

$$\begin{array}{c}
\frac{V \vdash e_1 \Downarrow v_1 \mid t_1 \quad V \vdash e_2 \Downarrow v_2 \mid t_2 \quad v_1, v_2 \in \mathbb{B}}{V \vdash e_1 \ \&\& \ e_2 \Downarrow v_1 \ \&\& \ v_2 \mid t_1 + t_2 + T_{BoolAnd}} \text{ (BOOLAND)} \\
\\
\frac{V \vdash e_1 \Downarrow v_1 \mid t_1 \quad V \vdash e_2 \Downarrow v_2 \mid t_2 \quad v_1, v_2 \in \mathbb{Z}}{V \vdash e_1 + e_2 \Downarrow v_1 + v_2 \mid t_1 + t_2 + T_{IntAdd}} \text{ (INTADD)} \\
\\
\frac{V \vdash e_1 \Downarrow v_1 \mid t_1 \quad \dots \quad V \vdash e_n \Downarrow v_n \mid t_n}{V \vdash (e_1, \dots, e_n) \Downarrow (v_1, \dots, v_n) \mid t_1 + \dots + t_n + T_{tupleHead} + nT_{tupleElem}} \text{ (TUPLE)} \\
\\
\frac{V \vdash tup \Downarrow (v_1, \dots, v_n) \mid t_1 \quad V[x_1 \mapsto v_1] \dots [x_n \mapsto v_n] \vdash e \Downarrow v \mid t_2}{V \vdash \text{let } (x_1, \dots, x_n) = tup \text{ in } e \Downarrow v \mid t_1 + t_2 + n T_{tupleMatch}} \text{ (TUPLEMATCH)} \\
\\
\frac{v = (V, \lambda x. e) \quad |FV(e) \setminus \{x\}| = n}{V \vdash \lambda x. e \Downarrow v \mid T_{funDef} + n T_{closure}} \text{ (CLOSURE)} \\
\\
\frac{V \vdash e_1 \Downarrow (V', \lambda x. e') \mid t_1 \quad V \vdash e_2 \Downarrow v_2 \mid t_2 \quad V'[x \mapsto v_2] \vdash e' \Downarrow v \mid t_3 \quad \mathbf{tag}(e_1) = \mathbf{tail}}{V; TM \vdash \text{app}(e_1, e_2) \Downarrow v \mid t_1 + t_2 + t_3 + T_{tailApp}} \text{ (TAILAPP)} \\
\\
\frac{V \vdash e_1 \Downarrow (V', \lambda x. e') \mid t_1 \quad V \vdash e_2 \Downarrow v_2 \mid t_2 \quad V'[x \mapsto v_2] \vdash e' \Downarrow v \mid t_3 \quad \mathbf{tag}(e_1) = \mathbf{normal}}{V; TM \vdash \text{app}(e_1, e_2) \Downarrow v \mid t_1 + t_2 + t_3 + T_{funApp}} \text{ (FUNAPP)}
\end{array}$$

Fig. 1. Selected rules of the big-step operational cost semantics.

with resource cost t . To differentiate between normal function calls and tail calls, we perform a semantics-preserving program transformation, which adds a tag to all function calls. A tag **tail** is added to tail calls and a tag **normal** is added to all other function calls. Our decision to give a positive cost to the constructs below, while a zero cost to other constructs comes from analyzing the compiled assembly code. Only the constructs below generated assembly instructions with a significant relative execution time. Intuitively, the cost of other constructs can be thought of as absorbed in these constructs. For example, the cost for addition absorbs the cost for loading the two addends.

Figure 1 contains illustrative example rules of the big-step cost semantics. The rules for operations on booleans and integers are very similar to **BOOLAND** and **INTADD**. For tuples, we introduce two constants $T_{tupleHead}$ and $T_{tupleElem}$. $T_{tupleHead}$ is counted every time we create a tuple, and $T_{tupleElem}$ is counted for the length of the tuple. Similarly, for tuple matching, we count a $T_{tupleMatch}$ for every element in the tuple being matched. When creating a tuple, there is an additional instruction, which assigns a tag to the tuple that represents the tuple constructor. Since there is no such instruction during a tuple match, we have an extra $T_{tupleHead}$ for creating tuples, but not when matching on it.

Since we support higher order functions, the rule **CLOSURE** for function definitions accounts for the cost of creating closures. We again introduce two constants

to deal with function definitions, T_{funDef} for creating a function, and $T_{closure}$ for creating a closure and capturing the free variables. Here, $FV(e)$ denotes the set of free variables of e . Since x is already bounded as the function argument, we remove it from the set of free variables of e . Rules `TAILAPP` and `FUNAPP` distinguish the cost of tail calls from the cost of regular function calls.

Lastly, we added a constant T_{base} to account for initializations made by each program, irrespective of the program code. Hence, we say that the execution cost of program p is $t + T_{base}$ if $\cdot \vdash p \Downarrow v \mid t$. With these evaluation rules for the cost semantics, we are ready to train our model by learning the values of the constructs described in this section.

4 Learning Memory Allocations

Before analyzing execution times, we will demonstrate the effectiveness of our approach by learning a cost semantics for memory allocations. We realized that our experiments did not lead to accurate results for floating point and tuple operations because they are generally stored on the heap, but often, optimized to be stored on the stack or the registers:

- The OCaml compiler performs constant propagation to determine which floats and tuples can be treated as globals, and need not be allocated on the heap, every time they are defined.
- If tuples only appear as arguments of a function, they are passed via registers and not allocated on the heap.

To accurately learn a cost semantics for floats and tuples we would need feedback from the OCaml compiler about the optimizations that have been performed. That is why, for the memory semantics only, we leave floats and tuples to future work and focus on booleans, integers, and lists. We use the cost semantics that is described in the previous section. According to this semantics, M , the number of bytes allocated by a program is a linear combination $M = \sum_{c \in \mathcal{C}} n_c M_c$.

We use the same training programs for learning memory allocations as for execution times. An interesting point is that the count n_c for each construct remains the same whether executing the training programs for time or memory. Hence, while performing the linear regression, we only need to execute the program on the interpreter once to obtain the counts n_c . We then use the `Gc` module in OCaml to obtain the number M of bytes allocated by the program. Since the memory allocation of a program is constant over different runs, we only need to measure the memory consumption once. For the Intel x86-64 platform, the memory costs of each construct obtained by the linear regression are as follows where $M_x = 0.00$ for all constants M_x that are not listed. $M_{base} = 96.03$ $M_{FunDef} = 24.00$ $M_{closure} = 7.99$ $M_{cons} = 24.00$

An analysis of the OCaml compiler indicates that rounding the learned constants to the nearest integer corresponds exactly to the number of bytes that are allocated by the corresponding construct. For example, our model implies that integers and booleans are not stored on the heap. And the OCaml manual [20]

indeed confirms that all integers and booleans are immediate, i.e., stored in registers and on the stack. The value 96 for the constant M_{base} is also confirmed as each program, even without memory allocations, has an initial heap consumption of 96 bytes. The cost $M_{FunDef} = 24$ and $M_{closure} = 8$ for function closures is also confirmed by the OCaml manual. If there are free variables trapped in the closure of a function, there is an additional memory allocation of 24 bytes on the heap to indicate that the program has a non-empty closure. Every cons constructor consumes 24 bytes on the heap; 8 bytes each for the head and tail, and 8 bytes for the tag to indicate the cons constructor. The empty list (`[]`) is represented as an immediate integer 0. Hence, the memory consumption of a list of size n is $24n$ bytes. Similarly, for the ARM v8-A platform, the memory costs of the non-zero constants obtained by the same linear regression are as follows. The results are also as expected and the data size seems to be 4 words.

$$M_{base} = 64.05 \quad M_{FunDef} = 12.00 \quad M_{closure} = 3.99 \quad M_{cons} = 12.00$$

We prefer learning the memory semantics instead of using them directly from the OCaml manual, because our technique is hardware-independent and can be extended in the event of creation of new architectures. It is notable that we can learn OCaml’s heap model by performing a simple regression without the constraint that the learned coefficients need to be integral or non-negative.

5 Learning Execution Times

As mentioned earlier, we used several regression techniques to train our cost semantics: linear regression, robust regression, and non-negative least squares. The accuracy of all three approaches is similar. Also, we train on the median execution times since they are less prone to noise than the mean. Below we give the cost of each high-level construct (in nanoseconds) trained using the normalized linear regression technique for the Intel x86-64 architecture. Intuitively, these constants define the median execution time of the respective construct on this specific platform.

$$\begin{array}{lll}
 T_{base} = 832.691 & T_{FunApp} = 1.505 & T_{TailApp} = 0.156 \\
 T_{FunDef} = 0.000 & T_{closure} = 2.921 & T_{BoolNot} = 0.424 \\
 T_{BoolAnd} = 0.184 & T_{BoolOr} = 0.183 & T_{IntUMinus} = 0.419 \\
 T_{IntAdd} = 0.297 & T_{IntSub} = 0.278 & T_{IntMult} = 1.299 \\
 T_{IntMod} = 19.231 & T_{IntDiv} = 19.011 & T_{FloatUMinus} = 1.232 \\
 T_{FloatAdd} = 2.102 & T_{FloatSub} = 2.116 & T_{FloatMult} = 1.737 \\
 T_{FloatDiv} = 8.575 & T_{IntCondEq} = 0.382 & T_{IntCondLT} = 0.381 \\
 T_{IntCondLE} = 0.381 & T_{IntCondGT} = 0.375 & T_{IntCondGE} = 0.381 \\
 T_{FloatCondEq} = 0.582 & T_{FloatCondLT} = 0.619 & T_{FloatCondLE} = 0.625 \\
 T_{FloatCondGT} = 0.585 & T_{FloatCondGE} = 0.629 & T_{letdata} = 2.828 \\
 T_{letlambda} = 1.312 & T_{letrec} = 1.312 & T_{patternMatch} = 0.223 \\
 T_{tupleHead} = 5.892 & T_{tupleElem} = 1.717 & T_{tupleMatch} = 0.237
 \end{array}$$

We make several qualitative observations about the learned cost semantics.

- $T_{FunApp} > T_{TailApp}$ indicates that a tail call is cheaper than a normal function call, confirming that tail call optimization reduces the execution time.
- $T_{BoolOr} \approx T_{BoolAnd}$, which is again expected as the $\&\&$ and $\|\$ operators just place jump instructions at appropriate locations in the assembly.
- $T_{IntMod} \approx T_{IntDiv} \gg T_{IntMult} > T_{IntAdd} \approx T_{IntSub}$. This is again expected, since whenever we perform integer division or modulo, a check is performed to see if the denominator is 0 and raise an appropriate exception. Hence, division and modulo are much more expensive than multiplication. The latter is more expensive than addition and subtraction.
- $T_{IntCondEq} \approx T_{IntCondLE} \approx T_{IntCondLT} \approx T_{IntCondGT} \approx T_{IntCondGE}$ is confirmed by studying the generated assembly code. A comparison is compiled to comparing the two integers and storing the result in a register, followed by a conditional jump. The analogous observation holds for floating point comparisons.

6 Garbage Collection

The OCaml garbage collector (GC) has 2 major components, the variable size major heap and the fixed size minor heap. Allocations occur first on the minor heap. If the minor heap is full, the GC is invoked. Roughly, the GC frees unreachable cells, and promotes all live variables on the minor heap to the major heap, essentially emptying the minor heap. OCaml employs a generational hypothesis, which states that young memory cells tend to die young, and old cells tend to stay around for longer than young ones.

We roughly model this behavior of the GC. Our hypothesis is that every call to the GC roughly starts with the full minor heap and empties the minor heap. We currently do not model the major heap. Hence, the time taken for each call to the GC in our model is roughly the same. We need two parameters to model the time taken by the GC, one is n_{gc} , which is the number of calls to the GC, and the other is T_{gc} , which is the time taken by 1 call to the GC. Our hypothesis states that $n_{gc} = \left\lfloor \frac{M}{H_0} \right\rfloor$, where H_0 is the size of the minor heap, and M is the total number of memory allocations. Since we can already model the number of heap allocations, all we need is the size of the minor heap.

Consequently, the two parameters T_{gc} and H_0 can be learnt from our training programs. OCaml offers a *Gc* module, which provides the number of calls to the GC. We use this module to find out the first call to the GC, the number of memory allocations and the time taken by the GC call, thereby learning H_0 , which is equal to the number of memory allocations (due to our hypothesis), and T_{gc} . With these parameters learned, the total execution time of the program is

$$T = \sum_{c \in \mathcal{C}} n_c T_c + \left\lfloor \frac{\sum_{c \in \mathcal{C}} n_c M_c}{H_0} \right\rfloor \cdot T_{gc}$$

7 Experiments with Control Programs

For our testing experiment, we used programs from the standard *List* library [2] and an OCaml tutorial [1].

Testing Method. For each test program, we first fix inputs of different sizes. For input i , we execute the compiled test program 500 times, measure the execution time, and compute the median T_i^{actual} . We then execute the test program on our interpreter to obtain the count of each construct n_c for $c \in \mathcal{C}$, and calculate $T_i^{expected} = \sum_{c \in \mathcal{C}} n_c T_c$. The average percentage error for each test program is

$$\text{Error}(\%) = \frac{1}{n} \left(\sum_{i=1}^n \frac{|T_i^{actual} - T_i^{expected}|}{T_i^{actual}} \right) \times 100$$

Experiments. Figure 2 shows the results on a compilation of control programs. The horizontal(x) axis represents the input size, while the vertical(y) axis represents the execution times. `factorialTR` is a tail recursive implementation of factorial. `append` concatenates two lists. `map` is a higher-order function, which maps a list of integers to booleans. Positive integers are mapped to `true` and the rest to `false`. `bubblesort` sorts a list using the bubblesort algorithm.

The measurement noise is significant, particularly in cases where high overhead is caused by context switches in the OS. Nevertheless, our prediction is surprisingly accurate given the simplicity of our model, particularly for functions without allocations, like `factorialTR`. For `append` our prediction is very accurate till the point of the first GC cycle ($x \approx 2.9$). The execution time after the first GC cycle is very unpredictable with frequent unexpected jumps (e.g. $x \approx 7$) and drops (e.g. $x \approx 14$). For `bubblesort` the GC jumps become invisible since the runtime of the GC is dominated by the actual computation. We can always very accurately predict the input size at which the GC cycles are triggered. This validates the accuracy of our model for learning memory allocations.

Table 1 summarizes the results obtained by evaluating our approach on 43 control programs. We implemented 3 different training algorithms for training in programs without GC, linear regression (LR), robust regression (RR) and non-negative least squares (NNLS). Each column represents the average percentage error for both architectures. We also tested all memory-intensive programs with larger inputs to evaluate our cost model for the GC. The last column presents the percentage error for programs with GC cycles, trained using linear regression. Note that the error increase of programs with GC cycles is not significant, and it indeed decreases for ARM architecture, indicating that our model for GC cycles is also accurate. However, in about 5% of the programs, the error is quite high (error above 50%). This is usually caused by an unmodeled compiler optimization, causing a significant reduction in execution time.

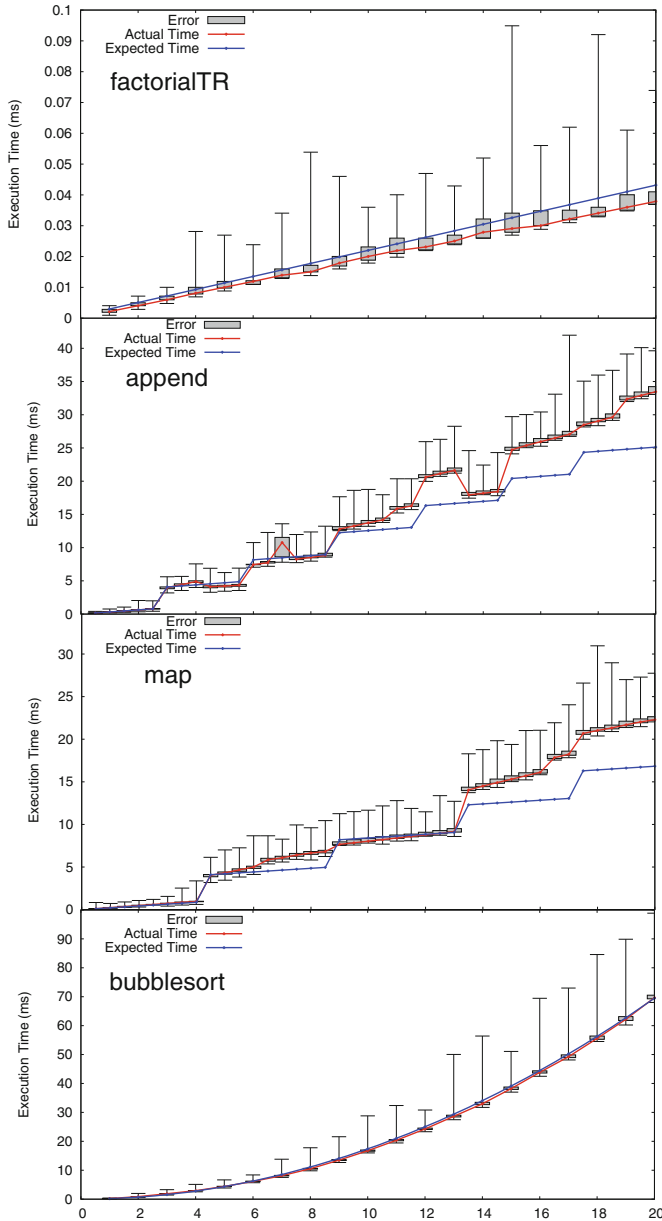


Fig. 2. Graph showing actual and expected time for factorialTR (input sizes $\times 10^3$) (top), append (input $\times 10^4$) (2nd), map (input $\times 10^4$) (3rd), and bubblesort (input $\times 10^2$) (bottom). The red and blue lines denote the actual and expected times, respectively. The vertical bars show the inherent noise in execution time. The lower and upper end of the vertical line denote the minimum and maximum execution time, while the lower and upper end of the box denotes the 1st and 3rd quartile of the execution time. (Color figure online)

Table 1. Results on x86-64 and ARM architectures

Architecture	Err (LR)	Err (RR)	Err (NNLS)	Err (GC)
x86-64	13.29	13.04	13.32	19.80
ARM v8-A	21.81	22.94	21.36	20.12

8 Applications

Integration with Resource Aware ML. We have integrated our learned cost semantics into Resource Aware ML [24], a static resource analysis tool for OCaml. RAML is based on an automatic amortized resource analysis (AARA) [23–25] that is parametric in a user defined cost metric. Such a metric can be defined by providing a constant cost (in floating point) for each syntactic form. This parametricity in a resource metric, and the ability to provide a cost to each syntactic form makes RAML very suitable for our integration purposes. Given an OCaml function, a resource metric, and a maximal degree of the search space of bounds, RAML statically derives a *multivariate resource polynomial* that is an upper bound on the resource usage as defined by the metric. The resource polynomial is parametric in the input sizes of the functions and contains concrete constant factors. The analysis is fully automatic and reduces bound inference to off-the-shelf LP solving. The subset of OCaml that is currently supported by RAML contains all language constructs that we consider in this paper. We used the experiments performed for this work to further refine the cost semantics and automatic analysis. For example, we added an analysis phase prior to the analysis that marks tail calls.

With the new cost metrics, we can use RAML for the first time to make predictions about the worst-case behavior of compiled code. For example, if we use the new execution-time metric (without GC) for x86 then we derive the following bounds in Table 2. The variables in the bounds are the input sizes. The table also contains runtime of the analysis and the number of generated constraint (cumulative for both bounds).

We can also use RAML to predict the execution time of programs with GC. To this end, we just derive two bounds using the execution metric and the metric for the number of allocations, respectively. We then combine the two bounds

Table 2. Symbolic bounds from RAML on x86-64

Program	Time bound (ns)	Heap bound (B)	Time (s)	#Cons
append	$0.45 + 11.28M$	$24M$	0.02	50
map	$0.60 + 13.16M$	$24M$	0.02	59
insertion sort	$0.45 + 6.06M + 5.83M^2$	$12M + 12M^2$	0.04	298
echelon	$0.60 + 17.29LM^2 + 23.11M + 37.38M^2$	$24LM^2 + 24M + 72M^2$	0.59	16297

using our model for GC which basically, accounts for a constant cost after a constant number of allocations. For example for `append` we obtain the following bound (heap size = 2097448 bytes, GC cycle = 3125429.15 ns on x86).

$$0.45 + 11.28M + \left\lfloor \frac{3125429.15 \times 24M}{2097448} \right\rfloor$$

Since the derived bounds for execution time and allocations are tight, this bound precisely corresponds to the prediction of our model as plotted in Fig. 2.

Qualitative Analysis. In addition to quantitative validation, we can also infer qualitative results from our learned cost semantics. For instance, we can compare two semantically-equivalent programs, and determine which one is more efficient on a specific hardware. Our model, predicts for example correctly the fastest version of different implementations of *factorial*, *append*, and *sieve of Eratosthenes*. Consider for example our Intel x86 machine and the following versions of `append`.

```

let rec append1 l1 l2 =
  match l1 with
  | [] -> l2
  | hd::tl -> hd::(append1 tl l2);

let rec append2 l1 l2 = match l1 with
  | [] -> l2
  | x::[] -> x::l2
  | x::y::[] -> x::y::l2
  | x::y::tl -> x::y::(append2 tl l2);

```

The trade-off in the above implementations is that the first has twice the number of function calls but half the number of pattern matches, as the second one. Since $T_{FunApp} = 1.505 > 4 \times 0.223 = 2 \times T_{patternMatch}$, hence, using our cost semantics concludes that the second version is more efficient. To reach this conclusion we can now analyze the two programs in RAML and automatically derive the execution-time bounds $0.45 + 11.28M$ and $0.45 + 10.53M$ for `append1` and `append2`, respectively. The fact that `append2` is faster carries over to the execution-time bounds with GC since the memory allocation bound for both functions is $24M$ bytes.

9 Related Work

The problem of modeling and execution time of programs has been extensively studied for several decades. Static bound analysis on the source level [4, 6, 9, 13, 15, 19, 22, 25, 35] does not take into account compilation and concrete hardware.

Closer related are analyses that target real-time systems by modeling and analyzing worst case execution times (WCET) of programs. Wilhelm et al. [37] provides an overview of these techniques, which can be classified into static [18], measurement-based methods, and simulation [8]. Lim et al. [29] associate a worst case timing abstraction containing detailed information of every execution path to get tighter WCET bounds. Colin and Puaut [17] study the effect of branch prediction on WCET. The goals of our work are different since we are not aiming at a sound bound of the worst-case but rather an approximation of the average case. Advantages of our approach include hardware independence, modeling of GC, and little manual effort after the cost semantics is defined.

Lambert and Power [27] introduced a hardware independent method of estimating time cost of Java bytecode instructions. Unlike our work, they do not take into account GC and compilation. Huang et al. [26] build accurate prediction models of program performance using program execution on sample inputs using sparse polynomial regression. The difference to our work is that they build a model for one specific program, are not interested in low-level features, and mainly want to predict the (high-level) execution time for a given input.

Acar et al. [3] learn cost models for execution time to determine whether tasks need to run sequentially or in parallel. They observe executions to learn the cost of one specific program. In contrast, we build a cost semantics to make predictions for all programs. There exist many works that build on high-level cost semantics [23], for instance to model cache and I/O effects [12]. However, these semantics do not incorporate concrete constants for specific hardware.

10 Conclusion and Future Work

We have presented an operational cost semantics learned using standard machine learning techniques, like linear regression, robust regression, etc. These semantics were able to model the execution time of programs with surprising accuracy; even in the presence of compilation and garbage collection. Since all the three models can be learned without relying on hardware specifics, our method is completely hardware independent and easily extensible to other hardware platforms. We have also presented an integration of the cost semantics with RAML, hence, allowing static analyzers to predict the execution time and heap allocations of assembly code for the first time.

One of the significant future directions is a more precise model for the garbage collector. Our model is limited to the minor heap, we need a model for the major heap and heap compactions as well. The size of the major heap is variable, hence, modeling the major heap is an important and complicated problem. We also need to incorporate other language features, especially user-defined data types in our semantics. Another challenge with real-world languages is compiler optimizations. We modeled one optimization (and suppressed another) in these semantics, but we should extend our semantics to incorporate more. Since these optimizations are performed at compile time, using static analysis techniques, it should be possible to model all of them. Finally, we think that it is possible to extend this technique to other programming languages, and we only need an appropriate interpreter to achieve that. We would like to validate this claim. We believe this connection between high-level program constructs and low-level program resources like time and memory is a first step towards connecting theoretical features of a language and its practical applications.

Acknowledgments. This article is based on research that has been supported, in part, by AFRL under DARPA STAC award FA8750-15-C-0082, by NSF under grant 1319671 (VeriQ), and by a Google Research Award. Any opinions, findings, and conclusions contained in this document are those of the authors and do not necessarily reflect the views of the sponsoring organizations.

References

1. 99 problems (solved) in OCaml. <https://ocaml.org/learn/tutorials/99problems.html>. Accessed 16 Aug 2016
2. Module list. <http://caml.inria.fr/pub/docs/manual-ocaml/libref/List.html>. Accessed 16 Aug 2016
3. Acar, U.A., Charguéraud, A., Rainey, M.: Oracle scheduling: controlling granularity in implicitly parallel languages. In: Proceedings of the 2011 ACM International Conference on Object Oriented Programming Systems Languages and Applications, OOPSLA 2011, pp. 499–518. ACM, New York (2011)
4. Albert, E., Arenas, P., Genaim, S., Gómez-Zamalloa, M., Puebla, G.: Automatic inference of resource consumption bounds. In: Bjørner, N., Voronkov, A. (eds.) LPAR 2012. LNCS, vol. 7180, pp. 1–11. Springer, Heidelberg (2012). doi:[10.1007/978-3-642-28717-6_1](https://doi.org/10.1007/978-3-642-28717-6_1)
5. Albert, E., Arenas, P., Genaim, S., Puebla, G., Zanardini, D.: Cost analysis of Java bytecode. In: Nicola, R. (ed.) ESOP 2007. LNCS, vol. 4421, pp. 157–172. Springer, Heidelberg (2007). doi:[10.1007/978-3-540-71316-6_12](https://doi.org/10.1007/978-3-540-71316-6_12)
6. Albert, E., Fernández, J.C., Román-Díez, G.: Non-cumulative resource analysis. In: Baier, C., Tinelli, C. (eds.) TACAS 2015. LNCS, vol. 9035, pp. 85–100. Springer, Heidelberg (2015). doi:[10.1007/978-3-662-46681-0_6](https://doi.org/10.1007/978-3-662-46681-0_6)
7. Alias, C., Darte, A., Feautrier, P., Gonnord, L.: Multi-dimensional rankings, program termination, and complexity bounds of flowchart programs. In: Cousot, R., Martel, M. (eds.) SAS 2010. LNCS, vol. 6337, pp. 117–133. Springer, Heidelberg (2010). doi:[10.1007/978-3-642-15769-1_8](https://doi.org/10.1007/978-3-642-15769-1_8)
8. Austin, T., Larson, E., Ernst, D.: SimpleScalar: an infrastructure for computer system modeling. *Computer* **35**(2), 59–67 (2002)
9. Avanzini, M., Lago, U.D., Moser, G.: Analysing the complexity of functional programs: higher-order meets first-order. In: 29th International Conference on Functional Programming (ICFP 2015) (2012)
10. Backes, M., Doychev, G., Köpf, B.: Preventing side-channel leaks in web traffic: a formal approach. In: Proceedings of 20th Network and Distributed Systems Security Symposium (NDSS) (2013)
11. Blanc, R., Henzinger, T.A., Hottelier, T., Kovács, L.: ABC: algebraic bound computation for loops. In: Clarke, E.M., Voronkov, A. (eds.) LPAR 2010. LNCS (LNAI), vol. 6355, pp. 103–118. Springer, Heidelberg (2010). doi:[10.1007/978-3-642-17511-4_7](https://doi.org/10.1007/978-3-642-17511-4_7)
12. Blelloch, G.E., Harper, R.: Cache and I/O efficient functional algorithms. In: Proceedings of the 40th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2013, pp. 39–50. ACM, New York (2013)
13. Brockschmidt, M., Emmes, F., Falke, S., Fuhs, C., Giesl, J.: Alternating runtime and size complexity analysis of integer programs. In: Ábrahám, E., Havelund, K. (eds.) TACAS 2014. LNCS, vol. 8413, pp. 140–155. Springer, Heidelberg (2014). doi:[10.1007/978-3-642-54862-8_10](https://doi.org/10.1007/978-3-642-54862-8_10)
14. Carbonneaux, Q., Hoffmann, J., Shao, Z.: Compositional certified resource bounds. In: 36th Conference on Programming Language Design and Implementation (PLDI 2015) (2015)
15. Černý, P., Henzinger, T.A., Kovács, L., Radhakrishna, A., Zwirchmayr, J.: Segment abstraction for worst-case execution time analysis. In: Vitek, J. (ed.) ESOP 2015. LNCS, vol. 9032, pp. 105–131. Springer, Heidelberg (2015). doi:[10.1007/978-3-662-46669-8_5](https://doi.org/10.1007/978-3-662-46669-8_5)

16. Chen, W.Y., Chang, P.P., Conte, T.M., Hwu, W.W.: The effect of code expanding optimizations on instruction cache design. *IEEE Trans. Comput.* **42**(9), 1045–1057 (1993)
17. Colin, A., Puaut, I.: Worst case execution time analysis for a processor with branch prediction. *Real-Time Syst.* **18**(2), 249–274 (2000)
18. Cousot, P., Cousot, R.: Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In: *Proceedings of the 4th ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages, POPL 1977*, pp. 238–252. ACM, New York (1977)
19. Danner, N., Licata, D.R., Ramyaa, R.: Denotational cost semantics for functional languages with inductive types. In: *29th International Conference on Functional Programming (ICFP 2015)* (2012)
20. Doligez, D., Frisch, A., Garrigue, J., Rémy, D., Vouillon, J.: The OCaml system release 4.03. <http://caml.inria.fr/pub/docs/manual-ocaml/>
21. Gulwani, S., Mehra, K.K., Chilimbi, T.M.: SPEED: precise and efficient static estimation of program computational complexity. In: *36th ACM Symposium on Principles of Programming Languages (POPL 2009)*, pp. 127–139 (2009)
22. Gulwani, S., Zuleger, F.: The reachability-bound problem. In: *Conference on Programming Language Design and Implementation (PLDI 2010)*, pp. 292–304 (2010)
23. Hoffmann, J., Aehlig, K., Hofmann, M.: Multivariate amortized resource analysis. In: *38th Symposium on Principles of Programming Languages (POPL 2011)* (2011)
24. Hoffmann, J., Das, A., Weng, S.: Towards automatic resource bound analysis for OCaml. In: *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages, POPL 2017, Paris, France, 18–20 January 2017*, pp. 359–373 (2017)
25. Hofmann, M., Jost, S.: Static prediction of heap space usage for first-order functional programs. In: *30th ACM Symposium on Principles of Programming Languages (POPL 2003)*, pp. 185–197 (2003)
26. Huang, L., Jia, J., Yu, B., Chun, B.G., Maniatis, P., Naik, M.: Predicting execution time of computer programs using sparse polynomial regression. In: Lafferty, J.D., Williams, C.K.I., Shawe-Taylor, J., Zemel, R.S., Culotta, A. (eds.) *Advances in Neural Information Processing Systems 23*, pp. 883–891. Curran Associates, Inc., Red Hook (2010)
27. Lambert, J.M., Power, J.F.: Platform independent timing of Java virtual machine bytecode instructions. *Electron. Notes Theor. Comput. Sci.* **220**, 97–113 (2008)
28. Lawson, C., Hanson, R.: *Solving Least Squares Problems*. Classics in Applied Mathematics. Society for Industrial and Applied Mathematics, Philadelphia (1995)
29. Lim, S.S., Bae, Y.H., Jang, G.T., Rhee, B.D., Min, S.L., Park, C.Y., Shin, H., Park, K., Moon, S.M., Kim, C.S.: An accurate worst case timing analysis for RISC processors. *IEEE Trans. Softw. Eng.* **21**(7), 593–604 (1995)
30. Neter, J., Kutner, M.H., Nachtsheim, C.J., Wasserman, W.: *Applied Linear Statistical Models*, vol. 4. Irwin, Chicago (1996)
31. Ngo, V.C., Dehesa-Azuara, M., Fredrikson, M., Hoffmann, J.: Quantifying and Preventing Side Channels with Substructural Type Systems (2016). Working paper
32. Olivo, O., Dillig, I., Lin, C.: Static detection of asymptotic performance bugs in collection traversals. In: *Conference on Programming Language Design and Implementation (PLDI 2015)*, pp. 369–378 (2015)
33. Rencher, A.C., Christensen, W.F.: *Multivariate regression*, pp. 339–383. Wiley (2012). <http://dx.doi.org/10.1002/9781118391686.ch10>
34. Rousseeuw, P.J., Leroy, A.M.: *Robust Regression and Outlier Detection*. Wiley, New York (1987)

35. Sinn, M., Zuleger, F., Veith, H.: A simple and scalable static analysis for bound analysis and amortized complexity analysis. In: Biere, A., Bloem, R. (eds.) CAV 2014. LNCS, vol. 8559, pp. 745–761. Springer, Heidelberg (2014). doi:[10.1007/978-3-319-08867-9_50](https://doi.org/10.1007/978-3-319-08867-9_50)
36. Steele Jr., G.L.: Debunking the ‘expensive procedure call’ myth or, procedure call implementations considered harmful or, LAMBDA: the ultimate GOTO. In: Proceedings of the 1977 Annual Conference, ACM 1977, pp. 153–162. ACM, New York (1977)
37. Wilhelm, R., et al.: The worst-case execution-time problem – overview of methods and survey of tools. *ACM Trans. Embed. Comput. Syst.* **7**(3), 36:1–36:53 (2008)
38. Zuleger, F., Gulwani, S., Sinn, M., Veith, H.: Bound analysis of imperative programs with the size-change abstraction. In: Yahav, E. (ed.) SAS 2011. LNCS, vol. 6887, pp. 280–297. Springer, Heidelberg (2011). doi:[10.1007/978-3-642-23702-7_22](https://doi.org/10.1007/978-3-642-23702-7_22)