

Performance Prediction upon Toolchain Migration in Model-Based Software

Aymen Ketata
Electrical and Computer Engineering
University of Waterloo
Waterloo, Canada
maketata@uwaterloo.ca

Carlos Moreno
Electrical and Computer Engineering
University of Waterloo
Waterloo, Canada
cmoreno@uwaterloo.ca

Sebastian Fischmeister
Electrical and Computer Engineering
University of Waterloo
Waterloo, Canada
sfischme@uwaterloo.ca

Jia Liang
Electrical and Computer Engineering
University of Waterloo
Waterloo, Canada
jliang@gsd.uwaterloo.ca

Krzysztof Czarnecki
Electrical and Computer Engineering
University of Waterloo
Waterloo, Canada
kczarnec@gsd.uwaterloo.ca

Abstract—Changing the development environment can have severe impacts on the system behavior such as the execution-time performance. Since it can be costly to migrate a software application, engineers would like to predict the performance parameters of the application under the new environment with as little effort as possible.

In this paper, we concentrate on model-driven development and provide a methodology to estimate the execution-time performance of application models under different toolchains. Our approach has low cost compared to the migration effort of an entire application. As part of the approach, we provide methods for characterizing model-driven applications, an algorithm for generating application-specific microbenchmarks, and results on using different methods for estimating the performance. In the work, we focus on SCADE as the development toolchain and use a Cruise Control and a Water Level application as case studies to confirm the technical feasibility and viability of our technique.

Index Terms—Model-based development, Migration, Automated Code Generation, Estimation, Prediction.

I. INTRODUCTION

Model Driven Development (MDD) is a software development approach where the source code can be automatically generated from the models. MDD is used to build an abstract representation of the system to improve productivity and communication among the managers, architects, designers, and developers. Some modeling languages such as UML, SysML, and modeling tools such as Eclipse modeling framework plugins are used in MDD development [22].

MDD requires specific toolchains to transform abstract models into executable code for the target system. A toolchain is a set of software tools used to generate code, compile and link, and it provides the binary code. In the context of real-time safety-critical applications, a qualified toolchain assures that the assertions and requirements claimed at the top-level are valid at the target programming language level.

Developers are reluctant to upgrade the toolchain while developing application models because, among other things,

even a small change in the toolchain can incur a migration cost and can add complexity in building the program. Changes in the toolchain can significantly affect the system behavior such as the execution performance of the program created with one toolchain. Various factors might motivate the migration decision to a new toolchain. For example, the new toolchain might offer new features that we want to benefit from or the old toolchain might lack support and maintenance. We focus in this paper on the performance parameters as the major decision factor of the migration. Predicting the performance parameters, before migrating the model to the new toolchain, can significantly help making the upgrade decision. Porting the application to a new toolchain should reflect a deep understanding of the performance changes under the new toolchain. The lack of structured and effective techniques to migrate the application to the new toolchain may lead to engineering work and an expensive migration with uncertain outcomes. This issue is critical in the embedded systems as the use of a new toolchain may also add safety and security issues. The upgrade of a verified toolchain such as the SCADE Systems is more difficult due to the strict requirements of the generated binary code [1], [18].

We provide a framework to predict the execution-time performance of a model-driven application on a new toolchain without migrating the entire application from the original toolchain. As the migration cost is in principle proportional to the amount of engineering work, our framework provides an automation process to extract and generate application-specific microbenchmarks. The execution-time estimates are relevant for the toolchain upgrade decision due to the trade-off between the cost and performance benefits. We present an automated technique to efficiently analyze the application model with respect to the new toolchain before the migration process as we can extrapolate and produce an accurate prediction of the execution-time performance.

The remainder of this paper proceeds as follows: Section II presents the problem statement and outlines our approach. Section III defines the used technical terms. In Section IV, we discuss several studies presenting ideas or techniques related to our work. Section V provides an overview of our developed tool. Finally, we describe our experimental setup in Section VI and results in Section VII. We discuss the validity and usability of our framework in Section VIII, followed by suggestions for future work and concluding remarks in sections IX and X.

II. OVERVIEW

This section introduces the problem that our work addresses and provides an overview of our proposed approach. Figure 1 illustrates the problem statement.

A. Problem Statement

Given two toolchains \mathcal{T}_1 and \mathcal{T}_2 and assuming a model has been developed and compiled to an executable under \mathcal{T}_1 , predict, with minimal porting effort expressed as the number of changes in the model, the execution-time performance, in the same execution environment, of the executable generated using \mathcal{T}_2 .

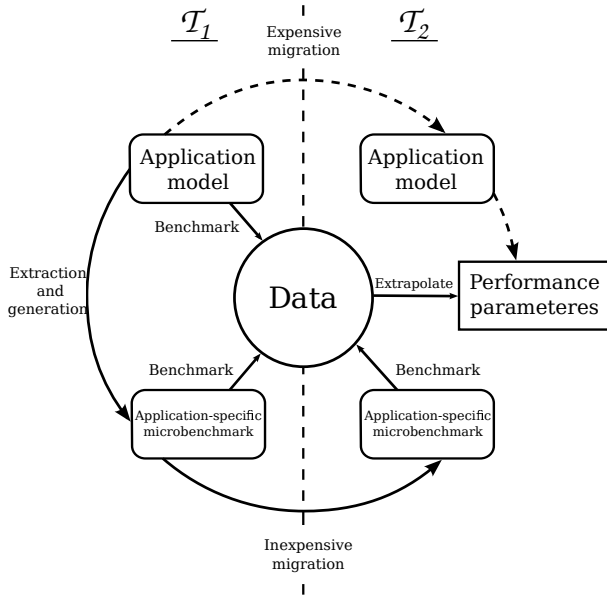


Fig. 1. Description of the problem statement

We refer to the solid lines path in Figure 1 as the anticlockwise path or inexpensive path. It represents the workflow of our approach to estimating the performance parameters. It is based on the extraction of metrics from the application model in \mathcal{T}_1 , and the generation of application-specific microbenchmarks in \mathcal{T}_1 . The next step is the migration of the microbenchmarks from \mathcal{T}_1 to \mathcal{T}_2 . We explain the details of this path in Section V. We refer to the dashed lines path as the clockwise path or expensive path. It represents the steps to follow if the entire application was migrated from \mathcal{T}_1 to \mathcal{T}_2 , which is in principle what we try to avoid for the purpose of performance prediction.

B. Our Approach

Our proposed approach allows developers to avoid porting the entire application when they face the decision of upgrading a development toolchain. Instead, we can predict, through benchmarking analysis, the performance of the application under the new toolchain \mathcal{T}_2 . This estimation is done at an early stage and avoids the cost of the migration of the entire application. The software artifacts such as the requirement, the specification, and the design documents that describe the abstraction of the final software system should be used to understand the application characteristics. To this end, we need to define the properties of the application under toolchain \mathcal{T}_1 that have an impact on the performance and derive application-specific microbenchmarks. To generate microbenchmarks that are representative of the application model, we use software metrics that capture the features and patterns in the application model that are relevant to its performance. Our approach focuses on the performance of model-based application relative to a toolchain rather than the application's architecture or design.

Our framework follows the inexpensive path from Figure 1. The framework is based on a set of tools that can automatically analyze the application, identify relevant characteristics in the model, and generate microbenchmarks that are representative of the original model. Users of the framework need to follow these steps:

- 1) **Measure application model with characteristic metrics:** The first step is to take the application model and characterize the model using a set of metrics. In our work we only use the ratio metric, however, other metrics such as Cyclomatic complexity, Halstead complexity, and Fan In/Fan-out [19], [27] are also applicable.
- 2) **Encode measurements in constraints:** The next step is to use the measurements and set up a series of constraints. These constraints will ensure that generated microbenchmarks are representative of the application model and maintain the original performance-relevant properties.
- 3) **Generate application-specific microbenchmarks:** The constraints permit multiple solutions. Depending on the required precision of the prediction, users then generate microbenchmark models from solutions to the constraint set. Being a statistical process, we can expect that increasing the number of microbenchmark models will lead to a lower standard error in the prediction
- 4) **Port microbenchmarks and generate code:** The migration of the microbenchmarks to the new toolchain should be relatively straightforward compared to the migration of the entire application model. Consequently, the porting effort will be inexpensive with respect to the required engineering effort. After porting, the user will compile the microbenchmarks on both toolchains and prepare the microbenchmarks for execution on the target platform.

- 5) **Benchmark microbenchmarks on the target platform:** By benchmarking the different models under the two toolchains and analyzing the results, the user can extrapolate an estimate of the execution-time performance of the entire application as migrated under the new toolchain.

Case studies of the SCADE toolchain for MDD provide evidence that our approach and framework are feasible. We used the Cruise Control and Water Level applications developed under two versions of SCADE systems and compared the predicted results to the migrated ones. Our estimates were reasonably accurate and correctly predicted that the executions of the applications under SCADE 6 were faster with respect to the applications under SCADE 5.

We remark that our approach does not predict the worst-case execution time (WCET). The presented method and tool constitute decision support to lower the risk when considering switching between toolchains. They focus on predicting average expected performance parameters while migrating toolchains and not on extremes like the WCET, which requires additional analysis and optimization. WCET analysis is required once the entire model is migrated to a different toolchain. For safety-critical systems, regulations most likely will require that a new analysis be made upon any change in the system, regardless of any estimates that would have been made prior to the changes.

III. BACKGROUND AND TERMINOLOGY

The SCADE Suite is a model-based development environment specifically tailored for safety-critical systems and often used in the avionics domain. The SCADE Suite is an integrated development environment that includes model-based design, simulation, verification, and qualified code generation. As SCADE provides a synchronous approach for reactive programs, it is suitable for developing safety-critical embedded software such as automotive and avionics applications [11].

Estimation or Parameter Estimation is the process of obtaining an approximate value of a parameter given available, and in general insufficient, data. In statistics, estimation usually refers to finding a function f of an observation vector X such that the error $|f(\mathbf{X}) - \theta|$ when estimating the parameter θ given an observation \mathbf{X} is minimized in some sense [28].

Information Extraction (IE) recognizes entities and relations from sets of data and transforms them into structured representations. In this paper, we focus on the extraction techniques used to generate application-specific microbenchmarks from the application models.

Migration is the conversion of a model developed under one version of a toolchain to either a new version of the same toolchain or a different toolchain. Unit, regression, and integration tests should be performed to validate the migration and ensure the required quality. The migration can affect the application behavior such as the performance parameters. The migration may be time-consuming, costly, and hard to perform

manually. The automation of the migration process could be used to reduce cost, but there are still important challenges to migrate concrete models [25].

Application Model refers to a set of large and complex models developed for a specific domain such as aircraft engines. Building such an application model requires the integration of physical, mathematical, and computational models.

An *application-specific microbenchmark* or *fingerprint model* is a reduced size model that shares common characteristics with the entire application model. Despite the fact that the microbenchmarks are randomly generated, they are representative of the application as they are constrained to a specified set of metrics extracted from the application.

Ratio metric is a software metric that consists of the fractions of each type of blocks in a model. That is, for each type of block, the number of blocks of that type divided by the total number of blocks is associated with the type. This metric does not consider I/O connections or the specifics of the design and structure of the model. Given a model, the computational cost of extraction of the ratio metric is linear with the number of blocks. We present a more detailed discussion and intuition on why this metric is relevant to the performance analysis in Section VI-B.

IV. RELATED WORK

The notion of software metrics is one of the key aspects in our method, as it is what captures the characteristics of a model that are relevant to its performance. Several studies exist in the literature that deal with this idea in the context of programming languages [7], [14], [15]. These software metrics include: Lines of Code, Cyclomatic complexity, Healdstead complexity, Cohesion and Coupling, Fan-In/Fan-Out and NPath. The Cyclomatic complexity [27] indicates the complexity of the application. It is computed based on the independent paths in the application generated with conditional statements. NPath complexity [10] measures the number of possible outcomes from the application. It might be hard to compute for large models that have nested conditional blocks. Fan-In/Fan-out [19] focuses on the information flow and measures the connections among the application components. Cohesion refers to the module responsibility and functionality as it expresses the degree of interdependency between the elements of the module. Coupling refers to the degree of dependency between the application modules. It explains the strength of the connection between the modules. Though these studies look for metrics that describe and capture the important characteristics of an application, like our work, they focus on program maintainability rather than performance and thus are not directly relevant to the problem that we are addressing.

We reviewed the applicability of software metrics to SCADE models. The work in [23] implemented a SCADE metric interpreter framework to extract the characteristics of the SCADE models. Some other related software metrics such as controllability and observability were presented in [12]. These metrics analyze the testability of the SCADE programs.

Testability metrics try to identify the different parts of the application that are critical, prone to errors, and difficult to validate. Similar ideas have also been investigated in the context of migrating legacy applications to modern programming languages with focus on the quality control of the application [24].

The survey [3] reviews research work that focuses on the performance prediction of model-based applications. The proposed approaches address the integration of performance analysis at early stages of the development process. Several works [16], [29] review model-based performance prediction approaches. These studies focus on the importance of conducting performance analysis at different stages throughout the software development cycle. Performance prediction requires a deep analysis of the system architecture from the requirement and specification phase to the configuration and deployment phase. Several works explore the performance requirements and try to predict the performance parameters at early stages to identify the issues of the concrete integration [13]. Some approaches, such as [20], focus on the application behavior after changes in source code and its impact on the application environment, dependencies, and performance. By contrast, our approach focuses on the changes introduced by the migration to a new toolchain.

Liu et al. [21] present an analytical approach that relies on stochastic modeling to predict the performance parameters of component-based applications. Our approach is statistical and based on measurement data from the actual target platform. The tool SoftArch/MTE [17] focuses on the evaluation of testbeds generated under various architectures of the application to help choose a particular architecture for the application design. In contrast, we estimate the effect of the migration between toolchains for the same application.

V. METHODS AND TOOLS

We now present the implementation of the framework and the methods and tools supporting it. To be able to evaluate the feasibility and viability of our concepts, we developed a tool that can generate application-specific microbenchmarks from a SCADE application model. Figure 2 gives an overview of the process.

We start by extracting the metric from the model. We focus on the ratio metric to analyze the models and identify its characteristics. Figure 3 is an example of a SCADE model. It has six math operators, one comparison operator, and one Boolean operator.

Based on the ratio metric, we developed a set of TCL scripts that analyze the application and extract the properties of the models. We encode the metamodels of the modeling language and the ratio metric extracted values as a set of constraints and feed them into a solver. We use Clafer [2] as the constraint language and its associated solver, which is based on the Choco constraint solver. Table I shows block occurrence and ratio data for the model in Figure 3.

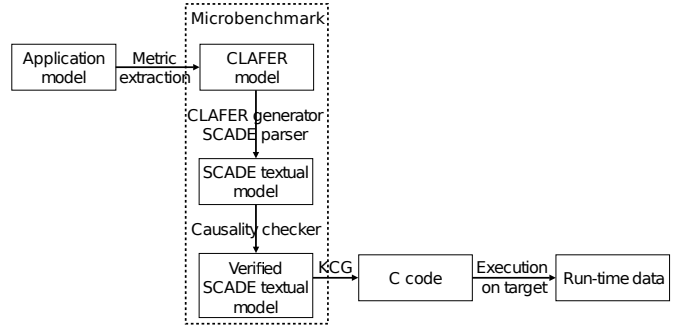


Fig. 2. Workflow generation

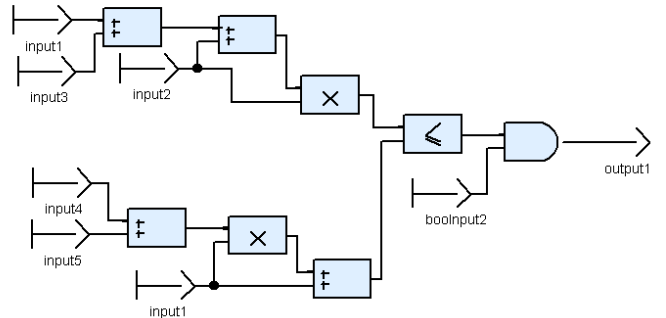


Fig. 3. SCADE model example

Given the extracted measurements, we produce constraints for the solver tool. For example, given that the size of the microbenchmarks is n , the number of Plus operators would be $50\% * n$, the number of Multiply operators would be $25\% * n$, and the number of And operators would be $12.5\% * n$

TABLE I: Extracted measurements data

	Occurrence	Ratio
Plus	4	50%
Multiply	2	25%
LessOrEqual	1	12.5%
And	1	12.5%

Clafer (**class, feature, reference**) is a lightweight modeling language with first-order relational logic. The Clafer compiler takes a model written in the Clafer modeling language as input and does some processing before invoking a backend solver to output instances conforming to the input model. Multiple solvers are supported such as SAT (Boolean satisfiability), SMT (satisfiability modulo theories), and CSP (constraint satisfaction problem). In this paper, the input is a SCADE metamodel written in Clafer, conjoined with the metric data constraints, and the output is a random SCADE model generated by the CSP solver.

Listing 1 shows an example of the Clafer representation of the constraint: we want to generate exactly one addition

operator either integer or real. Line 1 in Listing 1 specifies that `plus_Int` is an instance of the class `MathBlockInt`. The cardinality constraint `0..n` indicates that we would like to generate between 0 and `n` instances. Line 3 represents the constraint the sum of the number of instances of `plus_Int` and the number of instances of `plus_Real` must be equal to 1.

```
plus_Int : MathBlockInt 0..1
plus_Real : MathBlockReal 0..1
[#plus_Int + #plus_Real=1]
```

Listing 1. Metric data encoded in Clafer

Clafer generates random solutions, subject to the metric data and SCADE metamodel constraints, which we parse and convert to textual SCADE language. Textual SCADE is a declarative language; each line defines an element (e.g., one connection) in the model structure. Thus, the order of the lines is usually not relevant as in the case of imperative languages such as C.

For each set of constraints, we generated multiple random microbenchmarks. This is necessary to obtain a statistical characterization of the microbenchmark models' performance, which is necessary for the prediction. Listing 2 shows an example of the textual representation of one generated microbenchmark in the SCADE language. In this example, the first line indicates that we have an addition between the inputs `intInput1` and `intInput2`, and the result is stored in the variable `plus_Int0`.

```
function microbenchmark( /* Inputs , Outputs */
var
/* Local variables */
let
/* Update inputs */

plus_Int0 = intInput1 + intInput2 ;
plus_Int1 = plus_int0 + intInput3 ;
and_Bool0 = boolInput2 and LessOrEqual_Int0 ;
LessOrEqual_Int0 = multi_Int0 <= plus_Int0 ;
multi_Int0 = plus_Int1 * intInput4 ;
multi_Int1 = intInput5 * plus_Int0 ;
/* Update outputs */
tel
```

Listing 2. SCADE code of a sample microbenchmark

We generated C code using the toolchains KCG5.1 and KCG6.4 to compare the performance of the two toolchains. KCG is the automated code generator tool used to generate C code from SCADE models. KCG5.1 is the code generator tool used for SCADE 5 models, and KCG6.4 is the code generator tool used for SCADE 6 models. We automated the entire process flow to run the experiments. The scripts encode the metric, extract the application characteristics, generate SCADE microbenchmarks, run the SCADE checker, fix any causality errors by adding delay blocks, and generate C code. We compiled and executed the generated code for benchmarking. The DataMill infrastructure was used in our benchmarks to evaluate the performance of the C code generator of the SCADE toolchain. DataMill offers various architectures and software and hardware factors that can be used in the performance evaluation [9]. We used the i686 and x86_64

architectures to benchmark the C code, and we used the same hardware and software factors for the benchmarks.

A. The Constraint Solver

CSP [26] is a class of problems originating from the artificial intelligence community. A CSP problem is conventionally specified as a triple: V the set of variables, D the set of the variables' domains, C the set of constraints. A solution to a CSP problem is an assignment for each variable to a value in its domain such that none of the constraints are violated.

A CSP solver searches for solutions of a CSP problem by constructing an implicit search tree, where the variables are vertices and edges are assignments. The solver traverses the search tree in pre-order looking for leaf vertices such that every variable is assigned to a value, and none of the constraints are violated. These leaves are the solutions. The CSP backend for Clafer is implemented with the Choco library which supports integer variables and set variables over integers. Set variables are necessary and sufficient to encode the relational semantics of Clafer. For performance reasons however, the backend will optimize by using integer variables in place of set variables whenever possible. Choco's solving algorithm is based on an implicit search tree. The tree traversal can broadly be described in three steps starting at the root node:

- 1) Variable selection: Pick an unassigned variable using a heuristic. Most illustrations of CSP search trees would label the current node with the picked variable. Suppose the heuristic picked the integer variable i with domain $\{0, 1, 4\}$.
- 2) Decision: Assign the picked variable to a value in its domain and move down the current node's left branch. The right branch corresponds to the negative decision for when the algorithm *backtracks* to this node in the future. For example, if the left branch is $i = 1$ then the right branch is $i \neq 1$. More specifically, the left (respectively right) branch sets the domain of variable i to $\{1\}$ (respectively $\{0, 4\}$). There are other ways of making decisions such as domain splitting.
- 3) Constraint propagation: Infer new domains based on the available constraints. For example, if $i \neq j$ is a constraint, then remove 1 from the domain of variable j because assigning j to 1 will violate the constraint. If 1 is the only value in the domain of j , then the search entered a *contradiction* and can no longer proceed because the constraint $i \neq j$ is violated. The search then backtracks up the left branch(es) and goes down the nearest right branch and proceeds from there. Goto step 1 for the current node and repeat.

To generate random solutions, we modified the underlying CSP solver to construct the search tree randomly, i.e. the vertices/variables and edges/assignments are chosen randomly. Each search tree is used to generate only one solution. To generate n solutions, we generate n random search trees. The solutions from this approach are random in the sense that every

solution has a non-zero probability of being found. However, the probability distribution is not uniform.

B. Generation of Random Models

Model generation by searching random solutions to the constraints introduces some important challenges such as the potential production of invalid models. A model is considered invalid, if it does not meet the language requirements. In fact, even if the generated models satisfy the syntactic constraints of the modeling language, they might still not satisfy the semantic logic behind. In our study, SCADE is a synchronous language that guarantees that the data flow is immediately computed at each cycle with no physical latency. A loop in the generated model violates this rule, and the checker tool generates a causality error [4]. Such an error is not restricted to the use of SCADE and can occur in other synchronous modeling languages.

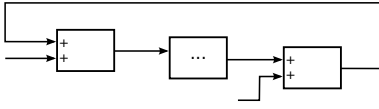


Fig. 4. Example model with causality error

In order to fix this issue, we developed a script that parses the generated model, and adds delay blocks to break the cycles [6].

SCADE imposes some restriction such as maintaining the same data flow and prohibiting mixed data type. Such requirements were encoded as a set of constraints in Clafer.

VI. EXPERIMENTATION

Given that we want to generate small-sized microbenchmarks that are easy to migrate, we fixed the size of the microbenchmarks. In that sense, the use of a constraint solver is not an issue.

A. Experimental Setup

The experiment was designed to run on DataMill. Several architectures are used for the benchmarking such as i686, and x86_64. We developed scripts to setup the experiment, run it and collect the data. The setup script uses GCC to compile the generated C code with no optimization flags. While running the experiment, we measured the execution-time of the microbenchmarks. The main function for the generated C code calls the main node of the microbenchmark. In other terms, the microbenchmark function would correspond to a call to one clock tick on the SCADE reduced model, and the data flow would pass through all the blocks of the model.

Figure 5 shows the execution-time measurement. In the first step, we warmed up the system to compensate for measurement errors due to memory caches, buffers, etc. This step is important as it warms up the caches in the system. In the next step, we measured the execution of the microbenchmark code several times. The last step is the calculation of the measured

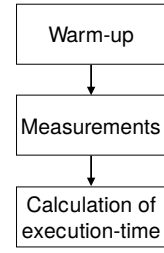


Fig. 5. Execution-time measurement

execution-time, and it is based on the different measurements of the previous step. The measured data is stored in log files. After collecting the data from DataMill, we performed sanity checks on the data. For example, we checked for non-negative execution-time values. We also inspected the distribution of the measured data by producing q-q plots and histograms. Figure 6 is an example that shows the q-q plot of the mean execution-times of 1000 microbenchmarks of the Water Level application. The x-axis shows the theoretical mean values, while the y-axis shows the sample values. The linearity of the points in the plot suggests that the data is normally distributed which is reasonable to expect.

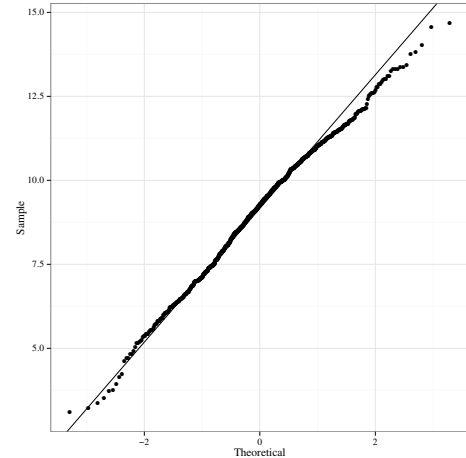


Fig. 6. Q-Q plot of the mean execution-time of the microbenchmarks of the Water Level application

As the input values of the model can affect the execution-time, the benchmarking process was designed to average out this effect. Random input values are generated for each execution of each microbenchmark model. Each microbenchmark is executed many times (10000 times) with randomly chosen inputs. The execution-time of one microbenchmark is given by the mean value of these executions. The use of random inputs allows the exploration of different paths in the code to obtain an “average behavior” that properly describes the average execution and average code coverage. If a probabilistic model of the inputs for the given application or application domain is available, such model should be used as the source of random values for the inputs in the benchmarking process. This has the advantage that the execution is now statistically representative

of the execution that the system will exhibit when in actual operation; thus, the prediction of the performance is specific to the real operating conditions under which the system will execute.

The hardware environment can also introduce bias in the execution performance. The execution-time can be mistakenly measured if executed on an exclusive architecture. The DataMill project solves this issue and proposes various architectures to benchmark the generated C code easily [9]. A set of scripts is used to setup the environment, run and collect the execution-time across several architectures.

B. Estimating Performance Parameters

Let \mathcal{P} be a model with $|\mathcal{P}|$ blocks, and consider toolchains \mathcal{T}_1 and \mathcal{T}_2 . Let b_1, b_2, \dots, b_c denote the C classes of blocks in the toolchains. For example, in SCADE, they would represent addition blocks, multiplication blocks, logical AND blocks, etc. Let $\tau_k^{(v)}$ be the execution-time of the fragment of code generated for blocks of type b_k under toolchain \mathcal{T}_v . We should expect execution-time for a particular type of block to be different under each toolchain, since the code generator is different. On the other hand, we assume a fixed execution-time for each type of block under the same toolchain, regardless of configuration and interaction with neighboring blocks. This is a reasonable approximation if the code generator is not highly optimizing, which is the case for tools that generate safety-critical qualified code such as SCADE.

Let p_k denote the fraction of blocks of type b_k in \mathcal{P} , and let $T^{(v)}$, with $v = 1, 2$ be the average execution-times of the model's main function generated by each of the toolchains \mathcal{T}_v . The averages are considered over the population of all possible models with $|\mathcal{P}|$ blocks that maintain the fractions p_k for each type of block. Notice that data-flow based models such as SCADE have a directed acyclic graph (DAG) representation [8]. Evaluation of the model's main function requires traversal of the graph, either in a breadth-first traversal until reaching all of the outputs or as a topological sort. In either case, the complexity of the operation is $O(V + E)$ where V is the number of vertices and E is the number of edges [8]. A further observation is that since blocks have inputs and outputs, and outputs cannot be "short-circuited" together, then each input can only be connected to one output. Since each block has $O(1)$ inputs, then $E = O(V)$, and thus evaluation of \mathcal{P} 's main function takes $O(|\mathcal{P}|)$ operations. Thus, we have

$$T^{(1)} = |\mathcal{P}| \sum_{k=1}^C \alpha_k^{(1)} p_k \tau_k^{(1)} \quad (1)$$

$$T^{(2)} = |\mathcal{P}| \sum_{k=1}^C \alpha_k^{(2)} p_k \tau_k^{(2)} \quad (2)$$

where the values $\alpha_k^{(v)}$ account for the average fraction of activity that each type of block in the model is exercised (including the multiplicative constant hidden in the big-Oh notation). If we consider a probability distribution of the

inputs that does not vary for the different models, then clearly the values $\alpha_k^{(v)}$ are fixed and determined by the various compatible ways to connect outputs of blocks to inputs of other blocks and the average paths of propagation of input data (which depend both on the structure of the model and the data).

Consider models with $|\mathcal{P}'|$ blocks that maintain the fractions of each type of block, p_1, p_2, \dots, p_c . It is reasonable to expect that the values of each $\alpha_k^{(v)}$ will be the same for models with $|\mathcal{P}'|$ blocks, since the fractions of each type of block p_k affect the possible configurations in which the sets of blocks can occur. Then, the average execution-times $T'^{(v)}$ for models of size $|\mathcal{P}'|$ are

$$T'^{(1)} = |\mathcal{P}'| \sum_{k=1}^C \alpha_k^{(1)} p_k \tau_k^{(1)} \quad (3)$$

$$T'^{(2)} = |\mathcal{P}'| \sum_{k=1}^C \alpha_k^{(2)} p_k \tau_k^{(2)} \quad (4)$$

From equations (1), (2), (3) and (4), we obtain

$$\begin{aligned} \frac{T^{(2)}}{T^{(1)}} &= \frac{T'^{(2)}}{T'^{(1)}} = \frac{\sum_{k=1}^C \alpha_k^{(2)} p_k \tau_k^{(2)}}{\sum_{k=1}^C \alpha_k^{(1)} p_k \tau_k^{(1)}} \\ \Rightarrow T^{(2)} &= T^{(1)} \frac{T'^{(2)}}{T'^{(1)}} \end{aligned} \quad (5)$$

If $\hat{T}^{(1)}$ and $\hat{T}^{(2)}$ are the execution-times for a given model (as opposed to the average execution-time over all possible models of this size), then Equation (5) yields an approximation, allowing us to obtain an estimate of the execution-time under \mathcal{T}_2 given a statistical representation of $T'^{(1)}$ and $T'^{(2)}$:

$$\hat{T}^{(2)} \approx \hat{T}^{(1)} \frac{T'^{(2)}}{T'^{(1)}} \quad (6)$$

Estimation of the variance is done in a similar way; the variances of $T^{(v)}$ are the result of the variances of the individual execution-times for the blocks, $\tau_k^{(v)}$, since all of the other terms are constants, provided that the the fractions of the types of blocks are preserved. Thus, assuming that the variables $\tau_k^{(v)}$ are uncorrelated, we have:

$$\text{Var} \left(T^{(v)} \right) = |\mathcal{P}| \sum_{k=1}^C \beta_k^{(v)} \text{Var} \left(\tau_k^{(v)} \right) \quad (7)$$

where the values $\beta_k^{(v)}$ depend exclusively on the values of $\alpha_k^{(v)}$ and the probability distribution of the $\tau_k^{(v)}$ variables.

Following a reasoning identical to that for the execution-time, we obtain a similar formula for the estimation of the variance of the execution-time under the new toolchain:

$$\text{Var} \left(\hat{T}^{(2)} \right) \approx \text{Var} \left(\hat{T}^{(1)} \right) \frac{\text{Var} \left(T'^{(2)} \right)}{\text{Var} \left(T'^{(1)} \right)} \quad (8)$$

1) *Estimating Ratios from \mathcal{T}_1 to \mathcal{T}_2* : Equations (6) and (8) give us an estimator for the parameters for models of one size based on the ratio of the parameters for models of a different size when migrated from \mathcal{T}_1 to \mathcal{T}_2 . Thus, we estimate these ratios based on sampling them through multiple randomly generated fingerprint models of a fixed size.

Each randomly generated model is executed multiple times with randomly selected input data to obtain an estimate of the execution-time’s mean and variance for the particular fingerprint model. We repeat this under both toolchains, to obtain one sample of the ratio between the parameters (mean and variance) under both toolchains. Assuming that both mean and variance, seen as random variables with respect to the population of all fingerprint models, follow a Gaussian distribution with non-zero mean, we empirically verified that their ratio follows a Gamma distribution. Thus, the mean of the samples obtained for each fingerprint model provides an adequate estimator for the required parameters.

VII. RESULTS

We present and discuss in this section the results of our case study benchmarks. We tried our approach on several applications provided as part of the SCADE software. We refer to these applications as validation models as they are provided as example applications, and the models are available in both version of SCADE. We used the SCADE example applications Cruise Control and Water Level. We developed another validation application that used a limited set of blocks, it used only Math operators such as addition and multiplication. For these applications, we had the models under SCADE 5 and the migrated models under SCADE 6. The microbenchmarks ratios represent the estimated results using our approach following the inexpensive path of Figure 1. The application ratios represent the ratios of the migrated models by following the expensive path of Figure 1. We calculate the mean execution-time of the different executions with random inputs. The application ratio is the mean execution-time of the entire application in SCADE 6 divided by the mean execution-time of the migrated application in SCADE 5. We compute the ratio of the execution-time for each microbenchmark. We refer to the mean value of the ratio of the microbenchmarks as the microbenchmarks ratio. We used the R boot package to compute the ratios and the 95% confidence intervals of the microbenchmarks [5].

Table II shows the mean execution-time benchmarking results of the validation models Math operators application, the Water Level application and the Cruise Control application. The rows show the target architectures that we used in our benchmarks, and the columns show the application ratio and the microbenchmarks ratio. The \pm refers to the margin of error computed by 95% confidence interval. These results are obtained by several executions of the C code with random inputs after a warm up phase of the system as explained in Section VI.

The Math operators application has the smallest difference between the estimated ratio (the microbenchmarks ratio) and the actual ratio value (the application ratio) in both architectures. This suggests that complex models may be subject to lower accuracy in the prediction. For example, the introduction of multiple execution paths by conditional blocks could justify such a difference. The benchmarking results suggest that the generated models were executed faster in SCADE 6 than SCADE 5. The entire application shows the same aspect for the evaluation of the execution-time. The above results show that the estimated ratios are consistent with the actual application ratios. Both the estimated and the migrated results show a speedup of the models when migrated to SCADE 6.

TABLE II: Execution-time ratio results
Architecture

	i686	x86_64
Math Operators		
Application Ratio	0.724±0.0007	0.664±0.002
Microbenchmarks Ratio	0.789±0.146	0.592±0.012
Water Level		
Application Ratio	0.489±0.001	0.454±0.001
Microbenchmarks Ratio	0.421±0.008	0.240±0.008
Cruise Control		
Application Ratio	0.791±0.011	0.325±0.001
Microbenchmarks Ratio	0.444±0.011	0.292±0.013

Even though the accuracy of the estimates is not too high, we notice that the framework gives a reasonable prediction of the performance evolution of the application. We notice that the 95% confidence intervals are tight which suggests that the effect of measurement errors and noise does not play a significant role in our results.

We can estimate that our method prediction is off by approximately 28% with respect to the correct performance for the migrated application (the geometric mean of actual ratio to predicted ratio is 1.28). A prediction accuracy of 28% can be useful in practice, since we can get severe performance regressions when migrating a model from one tool to another. Indeed, we observed performance variations between toolchains of over 400%; for example, for the microbenchmarks for the Water Level model, the ratio for the same microbenchmark executed under SCADE5 vs. executed under SCADE6 varied from approx. 0.2 to 0.9. Predicting performance change within 28% accuracy is a good starting point in the decision-making process. Furthermore, for all of the models that we used in our experiments, our technique correctly predicted whether performance would improve or deteriorate; we claim that trend prediction is as important or even more so than the exact percentage.

We omit the results for the variance prediction. We did extract these results from the experiments, but the figures that we obtained were unexpected and looked unreasonable. One potential reason for this is the completely different way in which SCADE 5 and SCADE 6 handle the conditional blocks;

from our observations of samples of generated code in both versions, conditional blocks in SCADE 5 involve execution of both branches followed by evaluation of the condition to choose one of the two already computed results. In SCADE 6, conditionals lead to optimized code: the condition is evaluated first, and only one of the two branches of the if-else is executed. This could have a profound impact on the variance of the execution-time, and, in particular, could have an effect on the accuracy of the prediction algorithm.

VIII. DISCUSSION

In this section, we discuss some of the important aspects that we believe have a critical impact on the validity and usability of our approach. Some of these issues could indeed represent threats to this validity. They could be related to the modeling language and the development toolchains. Below are some of the issues that we have considered:

- **Validity of software metrics.** The software metrics used to extract the data about the application models might not truly capture the application structure and architecture. For example, we might fail to detect the blocks patterns presented in the application. These patterns may affect the execution-time performance. Indeed, the interaction between connected blocks could have an important impact on the generated code's computational efficiency given low-level aspects such as cache, pipelining, or other hardware-related aspects. This makes us believe that there may be hard-to-capture underlying patterns in the structure that could have a significant impact on the overall performance. Thus, the accuracy of the prediction might benefit if the used metric would capture those patterns. The ratio metric proved to give reasonable results in our case studies, but we have to acknowledge the possibility that it might be insufficient for some other cases. We are convinced that this is the most important area that requires future work. Notice, however, that our proposed framework is extensible, and practitioners can use any other metrics that they have developed, and that may produce good results in the contexts being used.
- **Validity under different architectures.** We only benchmarked our experiments on a limited number of architectures. There are other architectures like ARM, MIPS, PowerPC, etc. However, Intel x86_64 and i686 cover an important fraction of the target audience for our method. Additional architectures will be used in future work.
- **Non-uniform random solutions.** The constraint solver does not generate perfectly random solutions, which could affect the accuracy of the prediction. We inspected the generated models and found a reasonable diversity in the models; we trust that this was not an issue in our experimentation. Moreover, generation of uniform random solutions by a constraint solver is a known complex problem, so our framework could certainly benefit from any progress that the AI community may make in this area.

- **Variance prediction sensitive to modeling tools.** As already mentioned in Section VII, the modeling tool could handle conditional and similar blocks in very different ways that could affect the prediction of the variance. This is a potentially critical aspect that we believe requires future work: changes in the variance may have an effect on estimates of WCET analysis if they were performed using measurement-based approaches. Even though this WCET analysis has to be done on the migrated model, the ability to accurately predict changes in the variance provides important information with respect to the risks that the migration could involve.
- **Restricted size of microbenchmarks.** Conceivably, the sizes that we chose for the microbenchmarks — which obey restrictions in the capacity of the constraint solver — could limit the accuracy of the predictions. The intuition is that larger microbenchmark models could have better ability to capture more complex characteristics. As constraint solvers become increasingly powerful, our framework could in turn benefit from any such advances.
- **Generalization.** We believe that our tool can be generalized and used with any modeling language. In fact, the benchmark generator uses a model similar to UML class models and OCL and thus can be used to model the complete syntax of any language. The grammar of any modeling language can be captured by our tool to generate syntactically valid models. As explained in Subsection V-B, the random generation may require the use of checking procedures to confirm the validity of the generated microbenchmarks. Our approach benefits from the nature of synchronous languages (such as limited use of iteration); the applicability of our approach to enterprise systems may not be obvious.

IX. FUTURE WORK

We aim to extend this work by studying the effect of other software metrics as mentioned in Section VI and the impact of the input distribution on the output one. We will study the impact of the data flow of the computations on the execution-time performance. This may be important as particular models in a given industrial application may have known distributions for the input data. A more detailed analysis of the effect of the distributions could have a positive impact in the applicability of our methodology. Furthermore, this could be useful when applying the technique to application domains instead of specific applications — for a given domain, typical distributions of input data may be known. We will also investigate the industrial practicality of our approach on turbojet engine applications with an avionics partner.

We think that our approach should be evaluated in other environments, such as Simulink. Comparing and combining our approach with analytical approaches are important future research directions.

X. CONCLUSIONS

We presented in this paper a framework to predict the execution-time performance parameters of model-based applications under different toolchains. Our approach has low cost as we avoid the migration cost by automating the migration of the models between the two toolchains. We can predict the performance parameters with minimal porting efforts expressed as the number of changes in the model.

To follow our framework, the application should be analyzed to extract software metrics that are relevant to the modeling language. The metric should be encoded in a constraint solver to generate application-specific microbenchmarks. The generated microbenchmarks are representative of the application, and the benchmarking of the microbenchmarks provides an estimate of the execution-time of the application under the two toolchains.

To illustrate our framework, we presented a SCADE Systems case study. We verified that our approach produced performance predictions that are reasonably close to the performance that we measure with concrete results.

XI. ACKNOWLEDGEMENTS

The authors would like to thank Michał Antkiewicz and Ed Zulkoski for fruitful discussions on constraint solvers and randomness for the solution, Yuguang Zhang and Jean-Christophe Petkovich for assistance with benchmarking on the DataMill platform.

REFERENCES

- [1] A. W. Appel. Verified Software Toolchain. In *Proceedings of the 20th European Conference on Programming Languages and Systems: Part of the Joint European Conferences on Theory and Practice of Software, ESOP'11/ETAPS'11*, pages 1–17, Berlin, Heidelberg, 2011. Springer-Verlag.
- [2] K. Bak, K. Czarniecki, and A. Wasowski. Feature and meta-models in clafar: mixed, specialized, and coupled. In *Software Language Engineering*, pages 102–122. Springer, 2011.
- [3] S. Balsamo, A. di Marco, P. Inverardi, and M. Simeoni. Model-based performance prediction in software development: a survey. *Software Engineering, IEEE Transactions on*, 30(5):295–310, May 2004.
- [4] G. Berry. SCADE: Synchronous Design and Validation of Embedded Control Software. In S. Ramesh and P. Sampath, editors, *Next Generation Design and Verification Methodologies for Distributed Embedded Control Systems*, pages 19–33. Springer Netherlands, 2007.
- [5] A. J. Canty. Resampling methods in r: the boot package. *R News*, 2(3):2–7, 2002.
- [6] J.-L. Colaço and M. Pouzet. Type-based Initialization Analysis of a Synchronous Data-flow Language. *Electronic Notes in Theoretical Computer Science*, 65(5):65 – 78, 2002. SLAP'2002, Synchronous Languages, Applications, and Programming (Satellite Event of {ETAPS} 2002).
- [7] S. D. Conte, H. E. Dunsmore, and Y. E. Shen. *Software Engineering Metrics and Models*. Benjamin-Cummings Publishing Co., Inc., Redwood City, CA, USA, 1986.
- [8] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to Algorithms, Third Edition*. The MIT Press, 2009.
- [9] A. B. de Oliveira, J.-C. Petkovich, T. Reidemeister, and S. Fischmeister. Datamill: Rigorous performance evaluation made easy. In *Proceedings of the 4th ACM/SPEC International Conference on Performance Engineering*, pages 137–148. ACM, 2013.
- [10] M. K. Debbarma, S. Debbarma, N. Debbarma, K. Chakma, and A. Jamatia. A Review and Analysis of Software Complexity Metrics in Structural Testing. *International Journal of Computer and Communication Engineering*, 2:129–133, 2013.
- [11] F.-X. Dormoy. SCADE 6: a model based solution for safety critical software development. In *Proceedings of the 4th European Congress on Embedded Real Time Software (ERTS'08)*, pages 1–9, 2008.
- [12] L. du Bousquet, M. Delaunay, H.-V. Do, and C. Robach. Analysis of testability metrics for Lustre/SCADE programs. In *Advances in System Testing and Validation Lifecycle (VALID), 2010 Second International Conference on*, pages 26–31. IEEE, 2010.
- [13] K. Falkner, V. Chiprianov, N. Falkner, C. Szabo, J. Hill, G. Puddy, D. Fraser, A. Johnston, M. Rieckmann, and A. Wallis. Model-Driven Performance Prediction of Distributed Real-Time Embedded Defense Systems. In *Engineering of Complex Computer Systems (ICECCS), 2013 18th International Conference on*, pages 155–158, July 2013.
- [14] N. E. Fenton and M. Neil. Software metrics: successes, failures and new directions. *Journal of Systems and Software*, 47(2–3):149–157, 1999.
- [15] N. E. Fenton and M. Neil. Software Metrics: Roadmap. In *Proceedings of the Conference on The Future of Software Engineering, ICSE '00*, pages 357–370, New York, NY, USA, 2000. ACM.
- [16] M. Fritzsche and J. Johannes. Putting Performance Engineering into Model-Driven Engineering: Model-Driven Performance Engineering. In H. Giese, editor, *Models in Software Engineering*, volume 5002 of *Lecture Notes in Computer Science*, pages 164–175. Springer Berlin Heidelberg, 2008.
- [17] J. Grundy, Y. Cai, and A. Liu. Softarch/mte: Generating distributed system test-beds from high-level software architecture descriptions. *Automated Software Engineering*, 12(1):5–39, 2005.
- [18] M. P. Heimdahl. Safety and software intensive systems: Challenges old and new. In *2007 Future of Software Engineering*, pages 137–152. IEEE Computer Society, 2007.
- [19] S. Henry and D. Kafura. Software structure metrics based on information flow. *Software Engineering, IEEE Transactions on*, (5):510–518, 1981.
- [20] R. Holmes and D. Notkin. Identifying program, test, and environmental changes that affect behaviour. In *Proceedings of the 33rd International Conference on Software Engineering*, pages 371–380. ACM, 2011.
- [21] Y. Liu, I. Gorton, and A. Fekete. Design-level performance prediction of component-based applications. *Software Engineering, IEEE Transactions on*, 31(11):928–941, 2005.
- [22] T. Mens and P. V. Gorp. A Taxonomy of Model Transformation. *Electronic Notes in Theoretical Computer Science*, 152(0):125 – 142, 2006. Proceedings of the International Workshop on Graph and Model Transformation (GraMoT 2005) Graph and Model Transformation 2005.
- [23] K. Ors and B. Laszlo. SCADE interpreter for measuring static and dynamic software metrics. In *Intelligent Systems and Informatics (SISY), 2013 IEEE 11th International Symposium on*, pages 123–128, Sept 2013.
- [24] G. Pandey, J. Jelschen, and A. Winter. Towards quality models in software migration. *analysis*, 1:M2.
- [25] L. M. Rose, M. Herrmannsdorfer, J. R. Williams, D. S. Kolovos, K. Garcés, R. F. Paige, and F. A. Polack. A comparison of model migration tools. In *Model Driven Engineering Languages and Systems*, pages 61–75. Springer, 2010.
- [26] F. Rossi, P. Van Beek, and T. Walsh. *Handbook of constraint programming*. Elsevier, 2006.
- [27] M. Shepperd and D. Ince. A critique of three metrics. *Journal of Systems and Software*, 26(3):197 – 210, 1994.
- [28] R. E. Walpole, R. H. Myers, S. L. Myers, and K. Ye. *Probability & Statistics for Engineers & Scientists*. Prentice-Hall, Ninth edition, 2011.
- [29] Q. Zhu and P. Deng. Design Synthesis and Optimization for Automotive Embedded Systems. In *Proceedings of the 2014 on International Symposium on Physical Design, ISPD '14*, pages 141–148, New York, NY, USA, 2014. ACM.