

# Tracing vs. Partial Evaluation

## Comparing Meta-Compilation Approaches for Self-Optimizing Interpreters

### Abstract

Tracing and partial evaluation have been proposed as meta-compilation techniques for interpreters. They promise that programs executing on simple interpreters can reach performance of the same order of magnitude as if they would be executed on state-of-the-art virtual machines with highly optimizing just-in-time compilers. Tracing and partial evaluation approach this meta-compilation from two ends of a spectrum, resulting in different sets of tradeoffs.

This study investigates both approaches in the context of self-optimizing interpreters, a technique for building fast abstract-syntax-tree interpreters. Based on RPython for tracing and Truffle for partial evaluation, we assess the two approaches by comparing the impact of various interpreter optimizations on the performance. The goal is to determine, whether either approach yields clear performance or engineering benefits. We find that tracing and partial evaluation both reach the same level of performance. With respect to the engineering, tracing has however significant benefits, because it requires language implementers to apply fewer optimizations to reach the same level of performance.

**Keywords** language implementation, just-in-time compilation, meta-tracing, partial evaluation, comparison, case study

### 1. Introduction

Interpretation is one of the simplest approaches to language implementation. However, interpreters lost some of their appeal because highly optimizing virtual machines (VMs) such as the Java Virtual Machine (JVM) or Common Language Runtime deliver performance that is multiple orders

of magnitude better. Nevertheless, interpreters still stand out for their simplicity, maintainability, and portability.

The development effort for highly optimizing static ahead-of-time or dynamic just-in-time compilers makes it often infeasible to build more than a simple interpreter. A recent example is JavaScript. In the last decade, its performance was improved by several orders of magnitudes, but it required major industrial investments. Unfortunately, such investments are rarely justified, especially for research projects or domain-specific languages (DSLs) with narrow use cases.

In recent years, tracing and partial evaluation became suitable meta-compilation techniques that alleviate the problem. RPython [Bolz et al. 2009; Bolz and Tratt 2013] and Truffle [Würthinger et al. 2012, 2013] are platforms for implementing languages based on simple interpreters that can reach the performance of state-of-the-art VMs. RPython uses trace-based just-in-time (JIT) compilation [Bala et al. 2000; Gal et al. 2006], while Truffle uses partial evaluation [Futamura 1971/1999] to guide the JIT compilation.

The PyPy<sup>1</sup> and Truffle/JS<sup>2</sup> projects show that general purpose languages can be implemented with good performance. However, for language implementers and implementation technology researchers, it remains the question of what the concrete tradeoffs between the two approaches are. When considering different purposes and maturity of language designs, the available engineering resources and the desired performance properties required different tradeoffs. For instance for a research language, it is most important to be able to experiment and change the language’s semantics. For the implementation of a standardized language however, the focus will typically be on performance and thus require the best possible mechanisms to realize optimizations. For implementation research, a good understanding of the tradeoffs between both approaches might lead to further improvements that simplify language implementation.

In this study, we compare tracing and partial evaluation as meta-compilation techniques for self-optimizing interpreters. We use RPython and Truffle as concrete representa-

<sup>1</sup>PyPy, a fast Python, access date: 2014-12-18 <http://pypy.org/>

<sup>2</sup>Truffle/JS, a JavaScript for the JVM, Oracle Labs, access date: 2014-12-18 <http://www.oracle.com/technetwork/oracle-labs/program-languages/javascript/index.html>

tions of these two approaches. To compare them in a meaningful way, we implement SOM [Haupt et al. 2010], a dynamic object-oriented language with closures, as identical as possible on top of both. Section 2 details the practical constraints and the requirements for a conclusive comparison. We investigate, which impact the two meta-compilation strategies have on a set of interpreter optimizations. The goal is to determine whether either of the two has clear advantages with respect to performance or engineering properties. The contributions of this paper are:

- a comparison of tracing and partial evaluation as meta-compilation techniques for self-optimizing interpreters.
- an assessment of the performance impact and implementation size of optimizations in self-optimizing interpreters.

We find that neither of the two approaches has a fundamental advantage for the reached peak-performance. However, meta-tracing has significant benefits from the engineering perspective. With tracing, the optimizer uses directly observed runtime information. In the case of partial evaluation on the other hand, it is up to the language implementer to capture much of the same information and expose it to the optimizer based on specializations.

## 2. Study Setup, Practical Constraints, and Background

The goal of this study is to compare tracing and partial evaluation as meta-compilation techniques with respect to the achievable performance as well as the required engineering effort for interpreters. This section gives a brief overview of meta-tracing and partial evaluation, and discusses how these two techniques can be compared based on concrete existing systems. It further discusses the design for the experimental setup, the concrete experiments, and the implications for the generalizability of the results. The section also provides the required background on self-optimizing interpreters and the SOM language, which we selected as case for this study.

### 2.1 Meta-Tracing and Partial Evaluation

While interpreters are a convenient and simple implementation technique, they are inherently slow. Hence, researchers tried to find ways to generate efficient native code from them without having to build custom JIT compilers. With the appearance of trace-based JIT compilation [Gal et al. 2006], trace-based meta-compilation, i. e., *meta-tracing* was the first practical solution for general interpreters [Bolz et al. 2009; Bolz and Tratt 2013]. The main idea is to trace the execution of the interpreter instead of tracing the concrete program it is executing, and thus, make the JIT compiler a reusable *meta-compiler* that can be used for different language implementations. The resulting traces are the units of compilation in such a system. Based on frequently executed loops on the application level, the interpreter records a con-

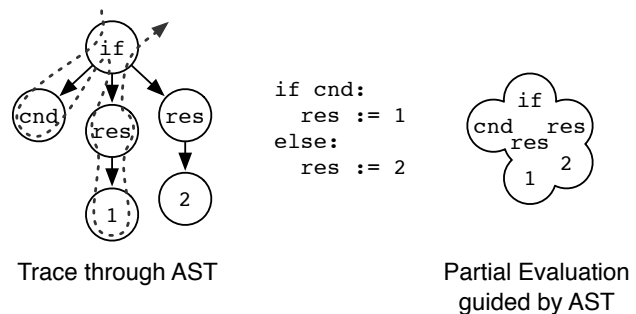
crete path through the program, which then can be heavily optimized and compiled to native code. Since traces span across many interpreter operations (cf. fig. 1), the interpreter overhead can typically be eliminated completely and only the relevant operations of the application program remain.

*Partial evaluation* [Futamura 1971/1999] of interpreters has been discussed as a potential meta-compilation technique for interpreters as well [Augustsson 1997; Sullivan 2001; Rigo and Pedroni 2006; Bolz et al. 2010]. However, only very recently, Würthinger et al. [2013] were able to show that it is a practical meta-compilation technique for abstract-syntax-tree-based (AST) interpreters. Instead of selecting the compilation unit by tracing, the unit is determined by using a program’s AST to guide a partial evaluator. The evaluator resolves all parts of the program that do not depend on unknown runtime information. With the knowledge of the AST and values embedded in it, the evaluator can resolve otherwise highly polymorphic method calls, perform aggressive constant propagation, and inlining. Thereby it identifies the relevant elements of the interpreter implementation (cf. fig. 1), which need to be included in a compilation unit.

In contrast to tracing, partial evaluation preserves the control flow of the interpreter and the user program that can not be resolved statically. Since the interpreter needs to handle every special case of a language, which leads to very complex control flow, partial evaluation and classic compiler optimizations alone were not able to generate efficient native code. Only with the idea of self-optimization, it became finally practical.

### 2.2 Self-Optimizing Interpreters

The main idea of a self-optimizing interpreter is that an executing AST rewrites itself at runtime, e.g., based on observed types and values [Würthinger et al. 2012]. Typical optimizations speculate for instance that observed types do not change in the future. For instance for an addition



**Figure 1.** Selecting JIT Compilation Units for AST Interpreters. To select a compilation unit, meta-tracing (left) records the operations performed by the interpreter for the execution of one specific path through a program. Partial evaluation (right) uses the AST structure to determine which interpreter-level code to include in a compilation unit.

operation, this allows to replace a generic node that handles all possible types by one specialized for integers. With such optimizations, an AST can specialize itself for exactly the way the program uses the language. This is beneficial for the interpreter, because it can avoid unnecessarily generic runtime operations, and at the same time the control flow is simplified, which leads to better compilation results when partial-evaluation-based meta-compilation is used [Würthinger et al. 2013].

Self-optimizations in general can also have other benefits. The previously mentioned type-based specialization of operations for instance avoids generic checks at runtime. Furthermore, it can be used to avoid boxing of primitive values such as integers to further reduce overhead and complexity of the operations. Another common optimization caches values for later use, e. g., with polymorphic inline caches for method lookups [Hölzle et al. 1991]. Starting out from a generic AST, the first execution of a method invocation node does the normal lookup and then rewrites itself to a simpler node that caches the lookup result and associates it with a predicate that confirms whether the cached value is valid in subsequent invocations. Thus, instead of having to include the complex lookup logic, the node only performs a check, and if it succeeds, the actual method invocation.

### 2.3 How to Compare Tracing and Partial Evaluation?

As discussed above, partial evaluation has only recently been shown to be practical and so far only in the context of self-optimizing interpreters. Meta-tracing has been successfully applied to AST interpreters as well [Bolz and Tratt 2013], thus, we compare both approaches based on self-optimizing AST interpreters.

To our knowledge RPython<sup>3</sup> is the only meta-tracing toolchain. Similarly, Truffle<sup>4</sup> is the only framework with partial-evaluation-based meta-compilation for interpreters. Thus, we chose these two systems for this experiment.

The goal of this study is to access the conceptual as well as the practical difference of tracing and partial evaluation. Hence, it stands to question what the generalizable insights of an empirical comparison are. From our perspective, both systems reached sufficient maturity and sophistication to represent the state of the art in tracing as well as partial evaluation technology. Furthermore, RPython with PyPy and Truffle with Truffle/JS implement complex widely used languages with the goal to optimize the peak performance as much as possible, and indeed reach the performance levels of dedicated JIT compiling VMs. Thus, we expect a performance comparison to reflect the general capabilities of the two approaches. However, both systems implement different sets of optimizations, and have different approaches for

generating native code. Therefore, minor performance difference between both systems are expected and will not allow for conclusions with respect to the general approaches. Nonetheless, we think the general order of magnitude will be representative for both approaches.

In order to compare both approaches fairly, we need a language implemented based on RPython as well as Truffle. With PyPy and ZipPy [Wimmer and Brunthaler 2013], there exist Python implementations for both systems. However, PyPy is a bytecode-interpreter and ZipPy a self-optimizing interpreter. Thus, a comparison would not only compare tracing with partial evaluation, but also include bytecode vs. ASTs, which would make a study inconclusive with respect to our question. The situation is the same for the Ruby implementations JRuby+Truffle<sup>5</sup> and Topaz. Moreover, they all differ in many other aspects, e. g., the set of implemented optimizations, which makes a comparison generally inconclusive. Hence, for a fair comparison we need language implementations for both systems that are as identical as possible, and enables us to compare tracing and partial evaluation instead of other aspects. For this study we use SOM, which is discussed in section 2.5.

### 2.4 RPython and Truffle

In the previous section, we discussed meta-tracing and partial evaluation from the conceptual perspective only. Since this study compares the two approaches empirically, this section provides a few technical details on RPython and Truffle.

*RPython* is a toolchain for language implementation that uses meta-tracing. It is also a restricted subset of Python and uses type inference and code transformations to add low-level services such as memory management and JIT compilation to interpreters to generate complete VMs. RPython's meta-tracing has been shown to work well for a wide range of different languages including Pyrolog (Prolog), Pycket (Racket), and Topaz (Ruby), of which some are bytecode interpreters, e. g., PyPy and Topaz, and others are AST interpreters, e. g., Pyrolog and Pycket.

With a set of annotations, language implementers can communicate high-level knowledge about the implemented language to the toolchain. Since trace-based compilation works best on loops, one of the main annotations is the so-called *trace merge point*, which indicates potential starting points for traces and defines how to recognize application-level loops. Other language-specific properties, for instance about mostly-constant values such as method lookup results can be communicated similarly. For instance, functions can have side-effects that are not essential for the execution, e. g., for caching the result of method lookups. With RPython's @elidable annotation, the optimizer can be told that it is safe to elide repeated executions within the context of a

<sup>3</sup>RPython Documentation, The PyPy Project, access date: 2015-03-18 <http://rpython.readthedocs.org/>

<sup>4</sup>The Truffle Language Implementation Framework, SSW JKU Linz, access date: 2015-03-18 <http://www.ssw.uni-linz.ac.at/Research/Projects/JVM/Truffle.html>

<sup>5</sup>JRuby+Truffle - a High-Performance Truffle Backend for JRuby, JRuby Project, access date: 2015-03-18 <https://github.com/jruby/jruby/wiki/Truffle>

trace. Another example are values that are runtime constants. Those can be explicitly *promoted* to enable the compiler to optimize based on them. In general, these annotations are useful in cases where an optimizer needs to make conservative assumptions, but the specific language usage patterns allow for optimistic optimizations, which can be used to generate specialized native code. A more detailed discussion of RPython is provided by [Bolz and Tratt \[2013\]](#).

*Truffle* is [Würthinger et al.](#)'s Java framework for self-optimizing interpreters and uses partial evaluation as meta-compilation technique. It integrates with the Graal JIT compiler for the partial evaluation of ASTs and the subsequent native code generation. Truffle in combination with Graal is built on top of the HotSpot JVM, and thus, guest languages benefit from the garbage collectors, memory model, thread support, as well as the general Java ecosystem.

For language implementers, Truffle provides an annotation-based DSL [[Humer et al. 2014](#)], which avoids much of the boilerplate code for typical self-optimizations. For instance, the DSL provides simple means to build specialized nodes for different argument types of operations. Instead of manually defining various node classes, the DSL provides a convenient way to provide only the actual operation. The node rewriting and argument checking logic is generated.

In addition to the DSL, there are other differences to RPython. For instance, runtime constants are exposed by providing node specializations instead of using a `promote`-like operation. Thus, the value is cached in the AST instead of relying on a trace context as RPython does. Another difference is that Truffle relies on explicit indications to determine the boundaries of compilation units. While RPython relies mostly on tracing, Truffle uses the `@TruffleBoundary` annotation to indicate that methods should not be included in the compilation unit. This is necessary, because Truffle uses a *greedy* inlining based on the partial evaluation, which would lead to too large compilation units without these explicit cutoffs. In practice, boundaries are typically placed on complex operations that are not on the fast path, e. g., lookup operations and complex library functionality such as string or hashtable operations. Also related is Truffle's `transferToInterpreter` operation, which results in a deoptimization point [[Hölzle et al. 1992](#)] in the native code. This excludes the code of that branch from compilation and can avoid the generation of excessive amounts of native code and enable optimizations, because the constraints of that branch do not have to be considered.

## 2.5 The Case Study: SOM (Simple Object Machine)

As discussed in [section 2.3](#), for a meaningful comparison of the meta-compilation approaches, we need close to identical language implementations on top of RPython and Truffle. We chose to implement the SOM language as case study. It is an object-oriented class-based language [[Haupt et al. 2010](#)] designed for teaching. Therefore, it is kept simple and in-

cludes only fundamental language concepts such as *objects*, *classes*, *closures*, and *non-local returns*. With these concepts, SOM represents a wide range of dynamic languages. Its implementation solves the same performance challenges more complex languages face, for instance for implementing exceptions, specializing object layouts, and avoiding the overhead for dynamic method invocation semantics, to name but a few.

While its size makes it a good candidate for this study, its low complexity raises the question of how generalizable the results of this study are for other languages. From our perspective, SOM represents the core concepts and thus solves many of the challenges common to more complex languages. What we do not investigate here is however the *scalability* of the meta-compilation approaches to more complex languages. Arguably, projects such as PyPy, Pycket, Topaz, JRuby+Truffle, and Truffle/JS demonstrate this scalability already. Furthermore, even so SOM is simple, it is a complete language. It supports classic object-oriented VM benchmarks such as DeltaBlue, Richards, and numeric benchmarks such as Mandelbrot set computation and n-body simulations. We further extended the benchmark set to include a JSON parser, a page rank algorithm, and a graph search to cover a wide range of use cases server programs might face.

**Implementation Differences of  $SOM_{MT}$  and  $SOM_{PE}$ .** Subsequently, we refer to the two SOM implementations as  $SOM_{MT}$  for the version with RPython's meta-tracing, and  $SOM_{PE}$  for one with Truffle's partial evaluation.  $SOM_{PE}$  builds on the Truffle framework with its TruffleDSL [[Humer et al. 2014](#)].  $SOM_{MT}$  however is built with ad hoc techniques to realize a self-optimizing interpreter, which are kept as comparable to  $SOM_{PE}$  as possible. Generally, the structure of the AST is the same for both interpreters. Language functionality such as method invocation, field access, or iteration constructs are represented in the same way as AST nodes.

Some aspects of the interpreters are different however.  $SOM_{PE}$  uses the TruffleDSL to implement basic operations such as arithmetics and comparisons. TruffleDSL significantly simplifies self-optimization based on the types observed at runtime and ensures that arithmetic operations can work directly on Java's primitive types `long` and `double` without requiring boxing. Boxing means that some primitive value is stored in a specifically allocated object. With Java's unboxed versions of primitive types, we avoid the additional allocation for the object and the pointer indirection when operating on the values.

$SOM_{MT}$  on RPython relies however on a uniform boxing of all primitive values as objects. With the absence of TruffleDSL for RPython, the minimal boxing approach used in  $SOM_{PE}$  was not practical because the RPython type system requires a common root type but does not support Java's implicit boxing of primitive types. Since tracing compilation eliminates the boxing within a compilation unit, it makes



only a difference in the interpreted execution and between compilation units. Thus, between compilation units, there might be additional overhead for allocations that are not done in  $SOM_{PE}$ . For the purpose of this paper, we consider this difference acceptable (cf. sections 4.3 and 4.4).

## 2.6 Assessing the Impact of the Meta-Compilation Strategies

To assess the benefits and drawbacks of meta-tracing and partial evaluation from the perspective of language implementers, we determine the impact of a number of interpreter optimizations on interpretation and peak performance. Furthermore, we assess the implementation sizes to gain an intuition of how the required engineering effort compares for both approaches.

**Optimizations.** To use a representative set of optimizations, we identify three main categories. *Structural optimizations* are applied based on information that can be derived at parse time. *Dynamic optimizations* require runtime knowledge to specialize execution based on observed values or types. *Lowerings* are reimplementations of standard library functionality that is performance critical inside the interpreter. These three groups cover a wide range of possible optimizations. For each category, we pick representative optimizations and detail them in section 3.

**Performance Evaluation.** For the performance evaluation, we consider the pure interpreted performance and the compiled peak performance. Both aspects can be important. While interpreter speed can be negligible for long-running server applications, it is critical for short-lived programs such as shell scripts. We assess the impact of the optimizations for both modes to also determine whether they are equally beneficial for interpretation and peak performance, or whether they might have a negative effect on one of them.

**Implementation Size of Optimizations.** To gain some *indication* for potential differences in engineering effort, we assess the implementation size of the applied optimizations. However, this is not a systematic study of the engineering effort. On the one hand RPython and Java are two very different languages making a proper comparison hard, and on the other hand, implementation size is only a weak predictor for effort. Nonetheless, implementation size gives an intuition and enables us to position the two approaches also with respect to the size of other language implementation projects. For instance in a research setting, an interpreter prototype might be implemented in 2.5K lines of code (LOC). A maturing interpreter might be 10 KLOC in size, but a state-of-the-art VM is usually larger than 100 KLOC.

## 3. Evaluated Optimization Techniques

The optimizations investigated in this study are chosen from the group of structural and dynamic optimizations as well as lowering of language and library functionality.

### 3.1 Structural Optimizations

Literature discusses many optimizations that can be performed after parsing a program, without requiring dynamic information. We chose a few to determine their impact in the context of meta-compilation. Note, each optimization has a shorthand by which we refer to it throughout the paper.

**Distinguish Variable Accesses in Local and Non-Local Lexical Scopes (*opt. local vars*)** In SOM, closures can capture variables of their lexical scope. A variable access thus needs to determine in which lexical scope the variable is to be found, then traverse the scope chain, and finally do the variable access. SOM's compiler can statically determine whether a variable access is in the local scope. At runtime, it might be faster to avoid the tests and branches of the generic variable access implementation. Thus, in addition to the generic AST node for variable access, this optimization introduces an AST node to access local variables directly.

**Handle Non-Local Returns Only in Methods including Them (*catch-return nodes*)** In recursive AST interpreters such as  $SOM_{PE}$  and  $SOM_{MT}$ , non-local returns are implemented using exceptions that unwind the stack until the method is found from that the non-local return needs to exit. A naive implementation handles the return exception simply in every method and checks whether it was the target. However, the setup for exception handlers as well as catching and checking an exception has a runtime cost on most platforms, and the handling is only necessary in methods that actually contain lexically embedded non-local returns. Thus, it might be beneficial to do the handling only in methods that need it. Since it is known after parsing a method whether it contains any non-local returns, the handling can be represented as an extra AST node that wraps the body of the method and is only added when necessary.

**Only Expose Variables in Lexical Scope that are Accessed (*min. escaping vars, SOM<sub>MT</sub> only*)** Truffle relies on a rigid framework that realizes temporary variables of methods with `Frame` objects. The partial evaluator checks that these *frames* do not escape the compilation unit, so that they do not need to be allocated. However, for lexical scoping, frame objects can *escape* as part of a closure object. In Truffle, such escaping frames need to be *materialized* explicitly. Instead of using such a strict approach, RPython works the other way around. An object can be marked as potentially *virtual*, so that its allocation is more likely to be avoided depending on its use in a trace.

To help the implicit approach of RPython in  $SOM_{MT}$ , the frames can be structured to minimize the elements that need to escape to realize closures. At method compilation time, it is determined which variables are accessed from an inner lexical scope and only those are kept in an array that can escape. The optimizer then ideally sees that the frame object itself does not need to be allocated. Since Truffle fixes the

structure of frames, this optimization is not applicable to  $SOM_{PE}$ .

***Avoid Letting Unused Lexical Scopes Escape (min. escaping closures)*** While the previous optimization tries to minimize the escaping of frames by separating variables, this optimization determines for the whole lexical scope whether it is needed in an inner scope or not. When the scope is not used, the frame object is not passed to the closure object and therefore will not escape. The optimization is realized by using a second AST node type that creates the closure object with `null` instead of the frame object.

### 3.2 Dynamic Optimizations

While the discussed static optimizations can also be applied to other types of interpreters, the dynamic optimizations are self-optimizations that require runtime information.

***Cache Lookup of Global Values (cache globals)*** In SOM, values that are globally accessible in the language are stored in a hash table. Since classes as well as for instance `true`, `false`, and `null` are globals, accessing the hash table is a common operation. To avoid the hash table lookup at runtime, globals are represented as association objects that can be cached after the initial lookup in a specialized AST node. The association object is necessary, because globals can be changed. For `true`, `false`, and `null`, we optimistically assume that they are not changed and specialize the access to a node that returns the corresponding constants directly.

***Caching Method Lookups and Block Invocations (inline caching)*** In dynamic languages, inline caching of method lookups is common to avoid the overhead of traversing the class hierarchy at runtime for each method invocation. In self-optimizing interpreters, this is represented as a chain of nodes, which encodes the lookup results for different kinds of objects as part of the caller’s AST. In addition to avoiding the lookup, this technique also exposes the target method as a constant to the compiler which in turn can decide to inline a method to enable further optimizations. Similar to caching method lookups, it is beneficial to cache the code of closures at their call sites to enable inlining.

In both cases, the node chains are structured in a way that each node checks whether its cached value applies to the current object or closure, and if that is not the case, it delegates to the next node in the chain. At the end of the chain, an uninitialized node either does the lookup operation or in case the chain grows too large, it is replaced by a fallback node that always performs the lookup.

***Type Specialization of Variable Accesses (typed vars,  $SOM_{PE}$  only)*** As mentioned earlier, Truffle uses `Frame` objects to implement local variables. For optimization, it tracks the types stored in a frame’s *slots*, i.e., of local variables. For  $SOM_{PE}$ , Truffle thus stores `long` and `double` values as primitives, which avoids the overhead of boxing. Furthermore,  $SOM_{PE}$ ’s variable access nodes specialize themselves

based on this type information to ensure that all operations in this part of an AST work directly with unboxed values.

Since  $SOM_{MT}$  uses uniform boxing, this optimization is not applied.

***Type Specialization of Argument Accesses (typed args,  $SOM_{PE}$  only)*** With the type specialization of  $SOM_{PE}$ ’s access to local variables, it might be beneficial to type specialize also the access to a method’s arguments. In Truffle, arguments to method invocations are passed as an `Object` array. Thus, this optimization takes the arguments passed in the object array and stores them into the frame object to enable type specialization. While this does not avoid the boxing of primitive values on method call boundaries, it ensures that they are unboxed and operations on these arguments are type specialized in the body of a method.

Note, since the variable access optimization is not applicable to  $SOM_{MT}$ , this optimization is not applicable either.

***Specialization of Object Field Access and Object Layout (typed fields)*** To avoid boxing, it is desirable to store unboxed values into object fields as well. Truffle provides support for a general object storage model [Wöß et al. 2014] that is optimized for class-less languages such as JavaScript, and is similar to maps in Self [Chambers et al. 1989]. To have identical strategies,  $SOM_{PE}$  and  $SOM_{MT}$  use a simplified solution that keeps track of how object fields are used at runtime, so that `long` and `double` values can be stored directly in the primitive slots of an object. For each SOM class, an object layout is maintained that maps the observed field types to either a storage slot for primitive values or to a slot for objects. The field access nodes in the AST specialize themselves according to the object layout that is determined at runtime to enable direct access to the corresponding slot.

***Specializing Array Representation (array strategies)*** Similar to other dynamic languages, SOM only has generic object arrays. To avoid the overhead of boxing, we implement strategies [Bolz et al. 2013] for arrays. It is similar to the idea of specializing the access and layout of object fields. However, here the goal is to avoid boxing for arrays that are used only for either `long`, `double`, or `boolean` values. In these cases, we specialize the storage to an array of the primitive type. In the case of booleans, it also reduces the size of the array from a 64-bit pointer to a byte per element.

***Specializing Basic Operations (inline basic ops.,  $SOM_{PE}$  only)*** As in other dynamic languages, SOM’s basic operations such as arithmetics and comparisons are normal method invocations on objects. Thus for instance the expression `1 + 2` causes the *plus* method to be invoked on the `1` object. While this allows developers to define for instance addition for arbitrary classes, in most cases arithmetics on numbers still use the built-in method. To avoid unnecessary method lookups and the overhead of method invocation, we specialize the AST nodes of basic operations directly to

the built-in semantics when the type information obtained at runtime indicate that it is possible.

Note, since this relies on TruffleDSL and its handling of the possible polymorphism for such specializations, this optimization is not applied to  $SOM_{MT}$ .

### 3.3 Lowerings

The last category of optimizations covers the reimplementa-tion of standard library functionality as part of the interpreter to gain performance.

**Control Structures (lower control structures)** Similar to specializing basic operations, we specialize control structures for conditionals and loops. In SOM, conditional structures are realized as polymorphic methods on boolean objects and loops are polymorphic methods on closures. An optimization of these constructs is of special interest because they are often used with lexically defined closures. Thus, in the context of one method, the closures reaching a control structure are statically known. Thus, specializing the control structures on the AST level does not only avoid overhead for method invocations done in the language-level implementation, but also utilizes directly the static knowledge about the program structure and exposes the closure code directly for further compiler optimizations such as inlining.

In  $SOM_{MT}$ , such specializations have the benefit of exposing the language-level loops to the implementation by communicating them directly to the meta-tracer with trace merge points (cf. section 2.4).

**Common Library Operations (lower common ops)** In addition to generic control structures, the SOM library provides many commonly used operations. We selected boolean, numeric, array copying, and array iteration operations for implementation at the interpreter level.

Similar to the specialization of basic operations and control structures, these optimizations are applied optimistically on the AST nodes that do the corresponding method invocation if the observed runtime types permit it.

## 4. Evaluation

Before discussing the results of the comparisons, we detail the methodology used to obtain and assess the performance and give a brief characterization of the used benchmarks.

### 4.1 Methodology

With the non-determinism in modern systems, JIT compilation, and garbage collection, we need to account for the influence of variables outside of our control. Thus, we execute each benchmark at least 500 times within the same VM instance. This guarantees that we have at least 100 continuous measurements for assessing steady state performance. The steady state is determined informally by examining plots of the measurements for each benchmark to confirm that the last 100 measurements do not show signs of compilation.

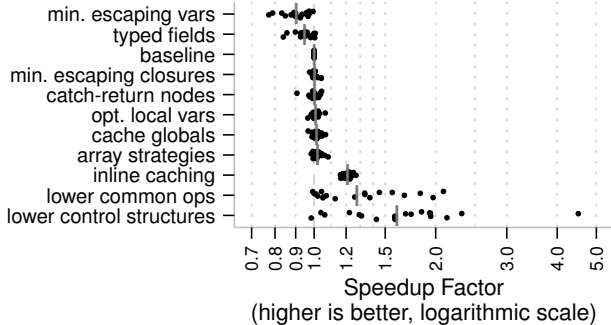
The benchmarks are executed on a system with two quad-core Intel Xeons E5520 processors at 2.26 GHz with 8 GB of memory and runs Ubuntu Linux with kernel 3.11, PyPy 2.4-dev, and Java 1.8.0\_11 with HotSpot 25.11-b03.

**Measurement Setup.** Pure interpretation performance for  $SOM_{MT}$  is measured with executables without meta-tracing support. Similarly, we measure the pure interpretation performance of  $SOM_{PE}$  on Hotspot without the partial evaluation and compilation support of Truffle. Thus, in both cases, there is no additional overhead, e. g., for compiler related bookkeeping. However,  $SOM_{PE}$  still benefits from the HotSpot’s normal Java JIT compilation, while  $SOM_{MT}$  is a simple interpreter executing directly without any underlying JIT compilation. We chose this setup to avoid measuring overhead from the meta-JIT compiler infrastructure and focus on the interpreter-related optimizations. Since we report results after warmup, the results for  $SOM_{PE}$  and  $SOM_{MT}$  represent the ideal interpreter performance in both cases.

For measuring the peak performance, we enable meta-compilation in both cases. Thus, the interpreter first executes with additional profiling overhead, and after completing a warmup phase, the benchmarks will execute solely in optimized native code. To assess the capability of the used meta-compilation approach, we report only the measurements after warmup is completed, i. e., ideal peak performance.

$SOM_{PE}$  uses a minimum heap size of 2GB to reduce noise from the GC. When measuring peak performance, Truffle is configured to avoid parallel compilation to be more comparable with RPython, which does not have any parallel execution. However, measurement errors for  $SOM_{PE}$  are generally higher than for  $SOM_{MT}$ , because the JVM performs various operations in parallel and reschedules the benchmark thread on other cores. RPython’s runtime system on the other hand is completely sequential, leading to lower measurement errors.

**Benchmark Suite.** The used benchmarks cover various aspects of VMs. DeltaBlue and Richards test among other things how well polymorphic method invocations are optimized. Json is a parser benchmark measuring string operations and object creation. PageRank and GraphSearch traverse large data structures of objects and arrays. Mandelbrot and n-body are classic numerical benchmarks focusing on floating point performance. Other benchmarks such as Fannkuch, n-queens, sieve of Eratosthenes, array permutations, bubble sort, and quick sort measure array access and logical operations. The storage benchmark is a stress test for garbage collectors. A few microbenchmarks test the performance, e. g., of loops, field access, and integer addition. While these benchmarks are comparably small and cannot compete with application benchmark suites such as Da-Capo [Blackburn et al. 2006], they test a relevant range of features and indicate the order of magnitude the discussed optimizations have on interpretation and peak performance.



**Figure 2.** Impact of optimizations on SOM<sub>MT</sub>'s interpreter performance. Experiments are ordered by geometric mean of the speedup over all benchmarks, compared to the baseline. Each dot represents a benchmark. The gray vertical bar indicates the geometric mean. The results show that the optimization for minimizing escaping variables slows the interpreter down. Inline caching and lowering of library functionality give substantial benefits.

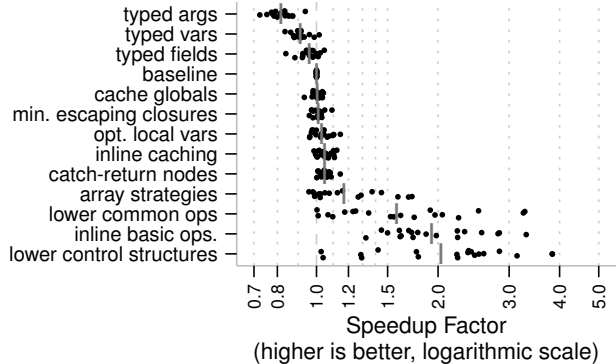
**Assessing Optimization Impact.** As in classic compilers, optimizations interact with each other, and varying the order in which they are applied can have significant implications on the observed gains they provide. To minimize the impact of these interdependencies, we assess the optimizations by comparing against a *baseline* that includes all optimizations. Thus, the obtained results indicate the gain of a specific optimization for the scenario where all the other optimizations have been applied already. While this might lead to underestimating the value of an optimization for gradually improving the performance of a system, we think it reflects more accurately the expected gains in optimized systems.

## 4.2 Impact on Interpreter

Before assessing the impact of the meta-compilation approach, we discuss the optimization's impact on interpretation performance.

Figure 2 depicts for each of the optimizations the benchmark results as separate points representing the average speedup over the baseline version of SOM<sub>MT</sub>. All dots on the right of the 1-line indicate speedup, while all dots left from the line indicate slowdowns. Furthermore, the optimizations are ordered by the geometric mean over all benchmarks, which is indicated for each optimization with a gray bar. Based on this ordering, all optimizations listed above the baseline cause on average a slowdown, while all optimizations listed below the baseline result in a speedup. Note, the x-axis uses a logarithmic scale.

The optimization for minimizing escaping of variables causes on average a slowdown of 9.6%. This is not surprising, since the interpreter has to allocate additional data structures for each method call and the optimization can only benefit the JIT compiler. Similarly, typed fields cause a slowdown of 5.3%. Since SOM<sub>MT</sub> uses uniform boxing, the inter-



**Figure 3.** SOM<sub>PE</sub> optimization impact on interpreter performance. Type-based specialization introduce overhead. Lowering of library functionality and direct inlining of basic operations on the AST-level are highly beneficial.

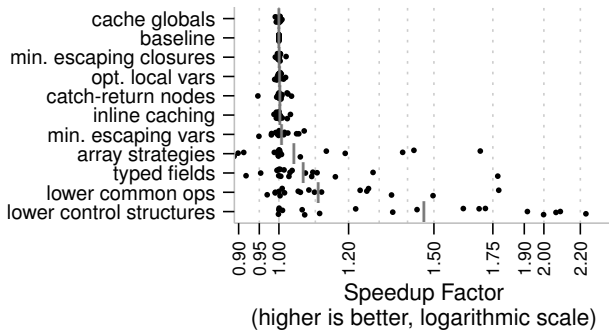
preter creates the object after reading from a field, and thus, the optimization is not beneficial. Instead, the added complexity of the type-specialization nodes causes a slowdown. The optimizations to separate catch-return nodes (0.2%), minimizing escaping of closures (0.2%), and the extra nodes for accessing local variables (0.8%) do not make a significant difference for the interpreter's performance. The dynamic optimizations for caching the association object of globals (1.4%) and array strategies (2%) do not provide a significant improvement either.

The remaining optimizations more clearly improve the interpreter performance of SOM<sub>MT</sub>. The largest gains for interpreter performance come from the lowering of control structures. Here we see an average gain of 1.6x (min. -1.6%, max. 4.5x). This is expected because their implementation in the standard library rely on polymorphic method invocations and the loop implementations all map onto the basic `while` loop in the interpreter. Especially for `for`-loops, the runtime overhead is much smaller when they are implemented directly in the interpreter because it avoids multiple method invocations and the counting is done in RPython instead of requiring language-level operations. Inline caching for methods and blocks (21%) gives also significant speedup based on runtime feedback.

For SOM<sub>PE</sub>, fig. 3 shows that the complexity introduced for the type-related specializations leads to overhead during interpretation. The typed arguments optimization makes the interpreter on average 18.3% slower. For typed variables, we see 8.9% overhead. Thus, if only interpreter speed is relevant, these optimizations are better left out. For typed object fields, the picture is less clear. On average, they cause a slowdown of 4.1%, but range from 16% slowdown to 4.5% speedup. The effect for SOM<sub>PE</sub> is more positive than for SOM<sub>MT</sub> because of the differences in boxing, but overall the optimization is not beneficial for interpreted execution.

Caching of globals (0.4%), optimizing access to local variables (3%), and inline caching (4.6%) give only minimal





**Figure 4.** SOM<sub>MT</sub> optimization impact on peak performance. Most optimizations do not affect average performance. Only lowering of library functionality gives substantial performance gains.

average speedups for the interpreter. The low gains from inline caching are somewhat surprising. However, SOM<sub>MT</sub> did not inline basic operations as SOM<sub>PE</sub> does. Thus, we assume that inlining of basic operations, which gives in itself a major speedup of 1.9x, hides the gains that inline caching of blocks and methods gives on an interpreter without it.

Array strategies give a speedup of 17.6% (min. -4.2%, max. 72.4%) and is with the different boxing strategy of SOM<sub>PE</sub> more beneficial for the interpreter. Similar to SOM<sub>MT</sub>, lowering library functionality to the interpreter level gives large improvements. Lowering common operations gives an average speedup of 1.6x and lowering control structures gives 2.1x, confirming the usefulness of these optimizations for interpreters in general.

### 4.3 Peak Performance

While some of the studied optimizations improve interpreted performance significantly, others cause slowdowns. However, especially the ones causing slowdowns are meant to improve peak performance for the meta-compilation with tracing or partial evaluation.

**Meta-Tracing.** Figure 4 shows the results for SOM<sub>MT</sub> with meta-tracing enabled. The first noticeable result is that 6 out of 10 optimizations have barely any effect on the optimized peak performance. The optimizations to cache globals (0%), minimize escaping closures (0.1%), optimize local variable access (0.2%), the separate nodes to catch returns (0.2%), inline caching (0.2%), and minimize escaping variables (0.7%) affect average performance only minimally.

For the optimization of local variable access and inline caching, this result is expected. The trace optimizer eliminate tests on compile-time constants and other unnecessary operations. Furthermore, inline caching is only useful for the interpreter, because SOM<sub>MT</sub> uses RPython’s `@elidable` (cf. section 2.4) to enable method lookup optimization. The lookup is marked as `@elidable` so that the optimizer knows its results can be considered runtime constants to avoid lookup overhead.

The optimization to minimize escaping of variables shows variability from a 5.1% slowdown to a to 6.8% speedup. Thus, there is some observable benefit, but overall it is not worth the added complexity, especially since the interpreter performance is significantly reduced.

Array strategies gives an average speedup of 4.7% (min. -29.9%, max. 69.3%). The additional complexity can have a negative impact, but also gives a significant speedup on benchmarks like the sorts that use integer arrays. For typed fields, the results are similar with an average speedup of 7% (min. -8.2%, max. 77.3%). For benchmarks that use object fields for integers and doubles, we see speedups, while others show small slowdowns from the added complexity.

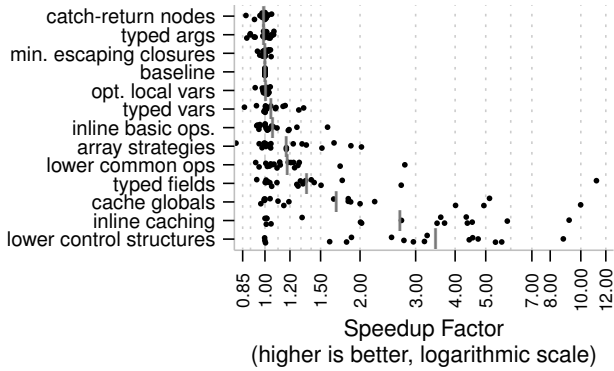
The lowering of library functionality is not only beneficial for the interpreter but also for meta-tracing. For common operations, we see a speedup of 11.5% (min. -21.6%, max. 1.8x). The lowering provides two main benefits. On the one hand, the intended functionality is expressed more directly in the recorded trace. For instance for simple comparisons this can make a significant difference, because instead of building, e. g., a *larger or equal* comparison with *smaller than* and negation, the direct comparison can be used. When layering abstractions on top of each other, these effects accumulate, especially since trace guards might prevent further optimizations. On the other hand, lowering typically reduce the number of operations that are in a trace and thus need to be optimized. Since RPython uses trace length as a criterion for compilation, lowering functionality from the library into the interpreter can increase the size of user programs that are acceptable for compilation.

For the lowering of control structures, we see a speedup of 1.5x (min. -0.1%, max. 4.1x). These speedups are based on the effects for common operations, but also on the additional trace merge points introduced for loop constructs. With these merge points, we communicate directly to RPython where user-level loops are and thereby provide more precise information for compilation.

Generally, we can conclude that only few optimizations have a significant positive impact when meta-tracing is used. Specifically, the lowering of library functionality into the interpreter helps to expose more details about the execution semantics, which enables better optimizations. The typing of fields and array strategies are useful, but highly specific to the language usage.

**Partial Evaluation.** The first observation based on fig. 5 is that compared to SOM<sub>MT</sub>, more of SOM<sub>PE</sub>’s optimizations have a positive effect on performance, which is also larger on average. Added catch-return node (-1.1%), typed arguments (-1.1%), minimization of escaping closures (-0.1%), and direct access to variables in local scope (0.3%) have only insignificant effect on peak performance.

Typed variables give an average speedup of (4.6%) (min. -13.9%, max. 32.6%). Thus, there is some speedup, how-



**Figure 5.**  $SOM_{PE}$  optimization impact on peak performance. Overall, the impact of optimizations in case of partial evaluation is larger. Lowering of control structures and inline caching are the most beneficial optimizations.

ever, in most situations partial evaluation is able to achieve the same effect without the type specialization.

Inlining of basic operations, which avoids full method calls, e.g., for arithmetic operations, shows a speedup of 5.8% (min.  $-5.8\%$ , max.  $1.6x$ ). It shows that in many cases the optimizer is able to remove the overhead of method calls. However, the optimization provides significant speedup in other cases as for instance complex loop conditions.

Array strategies give a speedup of 18.1% (min.  $-19\%$ , max.  $2x$ ), which is comparable to the speedup for  $SOM_{MT}$ , but slightly higher.

The lowering of common operations gives an average speedup of 18.7% (min.  $-6.5\%$ , max.  $2.8x$ ). The results are similar to the once for  $SOM_{MT}$ , indicating the general usefulness of these optimization independent of the technique to determine compilation units. Furthermore, the benefit of the optimization here is again higher for  $SOM_{PE}$ .

The optimization for object fields improves performance significantly. For the  $SOM_{PE}$  interpreter, it was causing a slowdown. With the partial evaluation and subsequent compilation however, we see a speedup of 41.1% (min.  $-5.8\%$ , max.  $11.2x$ ). Thus, typed object fields contribute significantly to the overall peak performance, despite their negative impact on interpreter performance. The benefit of typing variables and arguments seems to be however minimal. Here the optimizer has already sufficient information to generate efficient code regardlessly.

The caching of globals gives an average speedup of 79.9% (min.  $-3\%$ , max.  $10x$ ). Compared to RPython, on Truffle this form of node specialization is the only way to communicate runtime constants to the optimizer and as the results show, it is important for the overall performance.

Custom inline caching at method call sites and block invocations is the second most beneficial optimization. It results on average in a speedup of  $3x$  (min.  $0\%$ , max.  $19.6x$ ). On  $SOM_{MT}$ , this optimization did not give any improvements because RPython offers annotations that communi-

cate the same information to the compiler. With Truffle however, inline caching is only done by chaining nodes with the cached data to the call site AST node. While tracing intrinsically inlines across methods, Truffle needs these caching nodes to see candidates for inlining. Since inlining enables many other classic compiler optimizations, it is one of the the most beneficial optimizations for  $SOM_{PE}$ .

The lowering of control structures is the most beneficial optimization for  $SOM_{PE}$ . It gives an average speedup of  $4.3x$  (min.  $-0.2\%$ , max.  $232.6x$ ). Similar to  $SOM_{MT}$ , expressing the semantics of loops and other control flow structures results in significant performance improvements. In Truffle, similar to RPython, the control structures communicate additional information to the compilation backend. In  $SOM_{PE}$ , loops record loop counter to direct the adaptive compilation and branching constructs record branch profiles to enable optimizations based on branch probabilities.

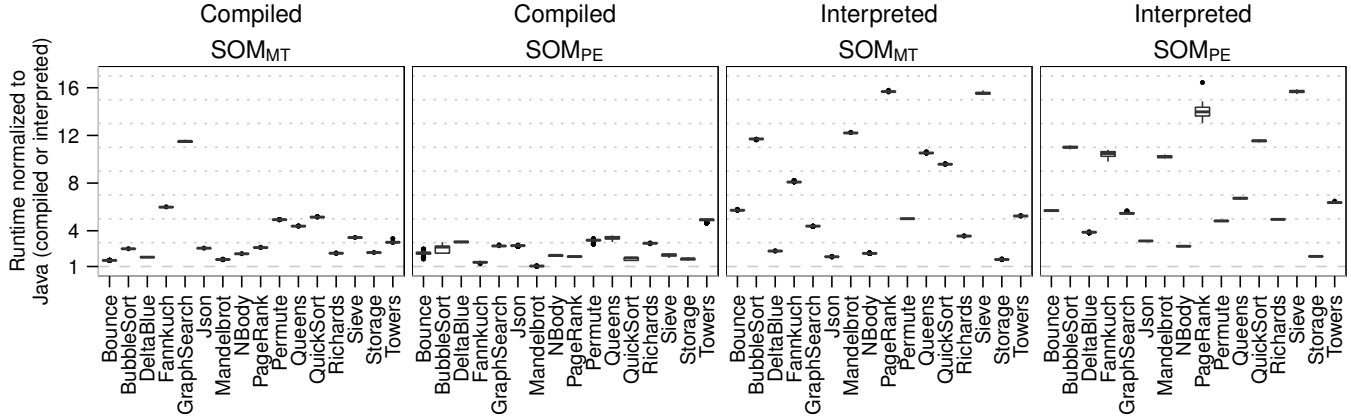
**Conclusion.** Considering all optimizations that are beneficial on average, and show for at least one benchmark larger gains, we find that array strategies, typed fields, and lowering of common operations and control structures are highly relevant for both meta-compilation approaches.

Inline caching and caching of globals is realized with annotations in RPython’s meta-tracing and thus, does not require the optimizations based on node specialization, even so, they are beneficial for the interpreted mode. However, with partial evaluation, the node specializations for these two optimizations provide significant speedup. Inlining of basic operations is beneficial for partial evaluation. While we did not apply this optimization to  $SOM_{MT}$ , it is unlikely that it provides benefits, since the same result is already achieved with the annotations that are used for basic inline caching. The typing of variables was also only applied to  $SOM_{PE}$ . Here it improves peak performance. For  $SOM_{MT}$ , it might in some cases also improve performance, but the added complexity might lead to a result like, e.g., for the minimizing of escaping variables, which does not improve peak performance on average.

Thus, overall we conclude that partial evaluation benefits more from the optimizations in our experiments by generating higher speedups. Furthermore, we conclude that more optimizations are beneficial, because partial evaluation cannot provide the same implicit specialization based on runtime information that meta-tracing provides implicitly.

#### 4.4 $SOM_{MT}$ vs. $SOM_{PE}$

To compare the overall performance of  $SOM_{MT}$  and  $SOM_{PE}$ , we use their respective baseline version, i.e., including all optimizations. Furthermore, we compare their performance to Java. The compiled performance is compared to the results for the HotSpot server compiler and the interpreted performance to the Java interpreter (`-Xint`). Note, the results for the compiled and interpreted modes are not comparable. Since the performance difference is at least one order



**Figure 6.** SOM performance compared to Java. The *compiled* performance are the SOMs with JIT compiler compared to HotSpot’s peak performance. The *interpreted* performance is compared to the HotSpot interpreter (`-Xint`).

of magnitude, the benchmarks were run with different parameters. Furthermore, cross-language benchmarking is inherently problematic. While the benchmarks are very similar, they are not identical and, the VMs executing them are tuned based on how the constructs are typically used, which differ between languages. Thus, the reported comparison to Java is merely an indication for the general order of magnitude one can expect, but no reliable predictor.

Figure 6 shows that  $SOM_{MT}$ ’s peak performance is on this benchmark set on average 3x (min. 1.5x, max. 11.5x) slower than Java 8 on HotSpot.  $SOM_{PE}$  is about 2.3x (min. 3.9%, max. 4.9x) slower. Thus, overall both SOMs reach the same order of magnitude of performance as Java, even so they are simple interpreters running on top of generic JIT compilation frameworks. Thus, the overall goal of achieving competitive performance is reached by both approaches. However,  $SOM_{MT}$  is slightly slower than  $SOM_{PE}$ . At this point, we are not able to attribute this performance difference to any conceptual differences between meta-tracing and partial evaluation as underlying technique. Instead, when investigating the performance differences, we see indications that the performance differences are more likely an indication of the amount of engineering that went into the RPython and Truffle projects, which results in Truffle and Graal producing more efficient machine code, while RPython has remaining optimization opportunities.

The performance of  $SOM_{MT}$  being only interpreted is about 5.6x (min. 1.6x, max. 15.7x) lower than that of the Java 8 interpreter. Similarly,  $SOM_{PE}$  is about 6.3x (min. 1.9x, max. 15.7x) slower than the Java 8 interpreter. Here we see some benchmarks being more than an order of magnitude slower. Such high overhead can become problematic when applications have short runtimes and very irregular behavior, because only parts of the application are executed as compiled code with good performance.

#### 4.5 Implementation Sizes

In addition to the achievable performance, engineering aspects can be of importance for language implementations as well. To gain some insight of how partial evaluation and meta-tracing compare in that regard, we determine the implementation sizes of the experiments. However, in addition to the weak insights measurement of implementation size provides, it needs to be noted that the obtained numbers are only directly comparable for experiments with the same SOM implementation. Since Java and RPython have significant syntactical and semantic differences, a direct comparison is not possible. Instead, we compare the relative numbers with respect to the corresponding baseline implementation. The reported percentages are based on the implementation without an optimization as denominator so that the percentage indicates the change needed to add the optimization.

As first indication, we compare the minimal versions of the SOM interpreters without optimizations with the baseline versions.  $SOM_{MT}$  has 3455 lines of code (LOC, excluding blank lines and comments) with all optimizations added it grows to 5414 LOC which is a 57% increase. The minimal version of  $SOM_{PE}$  has 5424 LOC and grows to 11037 LOC with all optimizations, which is an increase of 103%. Thus,  $SOM_{PE}$  is overall larger, which is expected since we apply more optimizations.

Table 1 lists the data for all experiments incl. absolute numbers. Comparing the relative increases of implementation sizes for  $SOM_{MT}$  and  $SOM_{PE}$  indicates that the optimizations are roughly of the same size in both cases. The only outlier is the implementation of inline caching which is larger for  $SOM_{PE}$ . Here the language differences between RPython and Java are becoming apparent and causes the  $SOM_{PE}$  implementation to be much more concise.

**Conclusion.** Considering performance and implementation sizes combined, we see for  $SOM_{MT}$  an overall peak

	SOM <sub>MT</sub> LOC %	SOM <sub>PE</sub> LOC %	SOM <sub>MT</sub>			SOM <sub>PE</sub>		
			LOC	ins.	del.	LOC	ins.	del.
baseline	0.0	0.0	5414	0	0	11037	0	0
typed args		1.4				10886	204	383
array strategies	11.6	9.0	4851	37	829	10125	126	1233
min. escaping closures	0.4	0.9	5394	5	30	10943	42	152
catch-return nodes	0.3	0.4	5397	12	36	10995	54	107
lower control structures	12.2	9.9	4824	8	790	10045	9	1160
inline caching	2.0	7.9	5307	1	158	10231	95	1095
inline basic ops.		3.7				10647	0	430
cache globals	0.5	1.7	5386	2	41	10853	14	239
opt. local vars	1.0	1.6	5359	49	135	10863	70	284
typed fields	10.2	11.1	4912	18	698	9933	39	1393
min. escaping vars	1.7		5322	20	130			
lower common ops	10.2	9.1	4912	2	678	10115	1	1083
typed vars		1.1				10915	9	161

**Table 1.** Implementation sizes of the implementations without the optimization. LOC: Lines of code excluding comments and empty lines, LOC %: increase of LOC to add optimization, ins./del.: inserted and deleted lines as reported by `git`

performance increase of 1.8x (min.  $-10.5\%$ , max. 5.4x) for going from the minimal to the baseline version. The interpreter performance improves by 2.4x (min. 41.5%, max. 3.9x). Note, the minimal version includes one trace merge point in the `while` loop to enable trace compilation (cf. section 2.4). For SOM<sub>PE</sub>, the peak performance improves by 78.1x (min. 22.8x, max. 342.4x) from the minimal to the baseline version. SOM<sub>PE</sub>’s interpreter speed improves by 4x (min. 2.1x, max. 7.3x). SOM<sub>PE</sub> also implements `while` in the interpreter, but it does not provide the same benefits for the partial evaluator as it does for the meta-tracer.

We conclude that for partial evaluation the optimizations are essential to gain performance. For meta-tracing however, they are much less essential and can be used more gradually to improve the performance for specific use cases.

## 5. Discussion

Section 2 already discussed general questions of the design of this study, its conclusiveness, and the generalizability of our results. Here we address a few more technical questions.

**Performance Results.** After studying the impact of various optimization on SOM<sub>MT</sub> and SOM<sub>PE</sub>, the question arises whether the observed performance effects are generalizable to other languages. Without further experiments, it needs to be assumed that they are not directly transferable. To give but a single example, for SOM<sub>PE</sub> we observed no benefit for peak performance from specializing method argument access based on their types. On the contrary, the interpreter showed clear performance drawbacks. However, in SOM, arguments are not assignable and methods are generally short. The usage pattern for arguments can thus be different in languages with assignable argument variables such as Java. Thus, other languages potentially benefit from this optimization. Nonetheless, the observations made here can

provide initial guidance for other language implementations to prioritize the optimization effort.

Another aspect that has not been studied is the impact of the optimizations on memory usage. The general requirement for self-optimizing interpreters is that the AST stabilizes at some point. This implies that self-modification should only introduce an upper bound of nodes, which limits the additional memory requirements. However, it has not yet been studied whether these optimizations cause significant additional memory consumption, which might be problematic in large applications.

**Meta-Tracing vs. Partial Evaluation.** A major difference between the two approaches is their overhead during interpretation. Partial evaluation requires the interpreter to record information about the executed code for instance branch probabilities and unused code paths. This information is used by the compiler to guide optimization together with heuristics, e. g., to avoid compilation of exception handling in the standard case. While sampling might reduce the overhead of collecting the runtime feedback, Truffle does currently use a precise approach that is active at all times, leading to a high overhead during interpretation.

With meta-tracing, the interpreter tracks execution only at the trace merge points. Only during tracing it records the additional information needed for optimization. Thus, interpretation performance might have conceptual advantages over a system with partial evaluation.

**RPython vs. Truffle.** When implementing a language, tooling can be a relevant deciding factor for RPython or Truffle. When optimizing an implementation, tools need to make it easy to understand and relate the optimizations done by the respective toolchains to an input program in the language that is implemented. Based on the current status of



the tools provided with both systems, there seems to be some benefit for meta-tracing. Since all optimizations are based on traces that linearize control flow, the tools are able to attribute relatively accurately the optimized instructions in a trace to the elements of the language implementation they originate from. In practice, this means that a program is relatively easily recognized in a trace, which supports the understandability of the results. For Truffle on the other hand, the available tool for inspecting the control- and data-flow graph of a program does not maintain the connection to the language implementation. Part of the issue is that some of Graal’s compiler optimizations can duplicate or merge nodes, which complicates the mapping to the input program.

Another practical aspect are the platforms’ capabilities and their ecosystems. Since Truffle builds on the JVM, support for threads, a memory model, and a wide range of software is implicitly given. Furthermore, the use of JVM-based software does not introduce a compilation boundary and thus, just-in-time compilation can optimize a Truffle-based language together with other libraries. RPython on the other hand does not come with comprehensive support for threads. Furthermore, its integration into the surrounding ecosystem is based on a foreign function interface (rffi), which is a compilation boundary for the tracing compiler.

One final difference between the two systems is their warmup performance. Currently, Truffle does not optimize the time it takes from starting an application until it reaches peak performance. Thus, we refrained from studying warmup performance. However, in our experience, RPython’s warmup is much faster than the warmup observed for Truffle. One aspect of that could be the extensive runtime profiling Truffle-based interpreters do. In [section 4.2](#), we reported the interpreter performance of SOM<sub>PE</sub> without this profiling enabled (cf. [section 4.1](#)), since it is currently not optimized. For compilation it is however essential and can slow down the interpreter significantly and thereby having a negative impact on overall application performance.

## 6. Related Work

As far as we are aware, there is no other study comparing meta-tracing and partial evaluation in detail. The closest related work is by [Marr et al. \[2014\]](#). They assess whether both approaches deliver on their promises, but, the comparison only regards the overall performance and compares in the case of meta-tracing a bytecode-based interpreter with a self-optimizing AST interpreter using partial evaluation.

Related to [Würthinger et al. \[2012\]](#)’s self-optimizing interpreters is for instance quickening and superinstructions focused on bytecode-based interpreters [[Proebsting 1995](#); [Casey et al. 2007](#); [Brunthaler 2010](#)].

The optimizations proposed for self-optimizing interpreters cover a wide range of topics and the optimizations discussed in this paper are either directly based on the literature or small variations. [Würthinger et al. \[2012\]](#) initially

discussed operation specialization by type, dynamic data type specialization, type specialization of local variable and field accesses, boxing elimination, and polymorphic inline caching (cf. also [Würthinger et al. \[2013\]](#)). Later, [Wöß et al. \[2014\]](#) detailed the strategy for field access optimization with an object storage model. [Kalibera et al. \[2014\]](#) discussed the challenges of a self-optimizing interpreter for the R language to address the dynamic and lazy nature of R. They detail a number of structural optimizations similar to the ones discussed here, dynamic operation and variable specialization, inline caching, data type specializations, as well as a profiling-based optimization of R’s view feature, which is a complex language feature that has different tradeoffs depending on the size of vectors it is used on. A similarly complex language feature that has been optimized in this context is Python’s generators [[Zhang et al. 2014](#)].

For meta-tracing, [Bolz and Tratt \[2013\]](#) discuss the impact on the VM design and implementation. They detail how an implementation needs to expose for instance data dependencies, compile time constants, and elidable computations clearly to the tracer for best optimization results. Generally, they advise to expose runtime constants also on the level of the used data structures. Thus, to prefer fixed sized arrays over variable sized lists, and to use known techniques such as *maps* [[Chambers et al. 1989](#)] to optimize objects to provide the tracer and subsequent optimization with as much information about runtime constants as possible. In this study, we find that these general suggestions apply to both compilation techniques, meta-tracing as well as partial evaluation.

## 7. Conclusion and Future Work

This study compares tracing and partial evaluation as meta-compilation techniques for self-optimizing AST interpreters. The results indicate that both enable language implementations to reach the same order of magnitude of performance as Java. A major difference between meta-tracing and partial evaluation is the amount of optimization a language implementer needs to apply to reach the same level of performance. Our experiments with SOM, a dynamic class-based language, indicates that meta-tracing performs well even without adding optimizations. With the additional optimizations it is on average only 3x (min. 1.5x, max. 11.5x) slower than Java. SOM<sub>MT</sub> reaches this result with 5414 LOC. For partial evaluation on the other hand, we find that many of the optimizations are essential to reach good performance. With all optimizations, SOM<sub>PE</sub> is on average only 2.3x (min. 3.9%, max. 4.9x) slower than Java. SOM<sub>PE</sub> reaches this result with 11037 LOC. Thus, we conclude overall from this study that meta-tracing and partial evaluation can reach the same level of performance. However, meta-tracing has significant benefits from the engineering perspective, because the optimizations provide generally fewer performance benefits and thus are less critical to be applied.

Since this study uses with Truffle and RPython two independent systems, we consider the observed difference in absolute performance as insignificant. We find that tracing and partial evaluation are equally suited for meta-compilation. The observed performance differences are merely an artifact of the different amounts of engineering that went into Truffle and RPython. Future work could verify this by studying both techniques on top of the same optimization infrastructure to determine whether for instance the preservation of control flow in compilation units is beneficial.

The interpreted performance of self-optimizing interpreters could still benefit from significant improvements. Possible research directions include approaches similar to superinstructions [Casey et al. 2007] on the AST level to avoid costly polymorphic method invocations. Another direction could be to attempt the generation of bytecode interpreters potentially in highly efficient machine code to reach interpretive performance competitive with for instance Java’s bytecode interpreter.

## Acknowledgments

Acknowledgments are left out in this double-blind version for review.

## References

- L. Augustsson. Partial Evaluation in Aircraft Crew Planning. In *Proc. of PEPM*, pages 127–136. ACM, 1997.
- V. Bala, E. Duesterwald, and S. Banerjia. Dynamo: A Transparent Dynamic Optimization System. In *Proc. of PLDI*, pages 1–12. ACM, 2000. ISBN 1-58113-199-2.
- S. M. Blackburn, R. Garner, C. Hoffmann, A. M. Khang, K. S. McKinley, R. Bentzur, A. Diwan, D. Feinberg, D. Frampton, S. Z. Guyer, M. Hirzel, A. Hosking, M. Jump, H. Lee, J. E. B. Moss, B. Moss, A. Phansalkar, D. Stefanović, T. VanDrunen, D. von Dincklage, and B. Wiedermann. The DaCapo Benchmarks: Java Benchmarking Development and Analysis. In *Proc. of OOPSLA*, pages 169–190. ACM, 2006. ISBN 1-59593-348-4.
- C. F. Bolz and L. Tratt. The Impact of Meta-Tracing on VM Design and Implementation. *Science of Computer Programming*, 2013.
- C. F. Bolz, A. Cuni, M. Fijalkowski, and A. Rigo. Tracing the Meta-level: PyPy’s Tracing JIT Compiler. In *Proc. of ICPOOLPS*, pages 18–25. ACM, 2009. ISBN 978-1-60558-541-3.
- C. F. Bolz, M. Leuschel, and D. Schneider. Towards a Jitting VM for Prolog Execution. In *Proc. of PDP*, pages 99–108. ACM, 2010. ISBN 978-1-4503-0132-9.
- C. F. Bolz, L. Diekmann, and L. Tratt. Storage Strategies for Collections in Dynamically Typed Languages. In *Proc. of OOPSLA*, pages 167–182. ACM, 2013.
- S. Brunthaler. Efficient Interpretation Using Quickening. In *Proc. of DLS*, pages 1–14. ACM, Oct. 2010.
- K. Casey, M. A. Ertl, and D. Gregg. Optimizing Indirect Branch Prediction Accuracy in Virtual Machine Interpreters. *ACM Trans. Program. Lang. Syst.*, 29(6):37, 2007. ISSN 0164-0925.
- C. Chambers, D. Ungar, and E. Lee. An Efficient Implementation of SELF a Dynamically-Typed Object-Oriented Language Based on Prototypes. In *Proc. of OOPSLA*, pages 49–70. ACM, 1989. ISBN 0-89791-333-7.
- Y. Futamura. Partial Evaluation of Computation Process—An Approach to a Compiler-Compiler. *Higher-Order and Symbolic Computation*, 12(4):381–391, 1971/1999. ISSN 1388-3690.
- A. Gal, C. W. Probst, and M. Franz. HotpathVM: An Effective JIT Compiler for Resource-constrained Devices. In *Proc. of VEE*, pages 144–153. ACM, 2006. ISBN 1-59593-332-6.
- M. Haupt, R. Hirschfeld, T. Pape, G. Gabrysiak, S. Marr, A. Bergmann, A. Heise, M. Kleine, and R. Krahn. The SOM Family: Virtual Machines for Teaching and Research. In *Proc. of ITiCSE*, pages 18–22. ACM Press, June 2010.
- U. Hölzle, C. Chambers, and D. Ungar. Debugging Optimized Code with Dynamic Deoptimization. In *Proc. of PLDI*, pages 32–43. ACM, 1992. ISBN 0-89791-475-9.
- C. Humer, C. Wimmer, C. Wirth, A. Wöß, and T. Würthinger. A Domain-Specific Language for Building Self-Optimizing AST Interpreters. In *Proc. of GPCE*, pages 123–132. ACM, 2014.
- U. Hölzle, C. Chambers, and D. Ungar. Optimizing Dynamically-Typed Object-Oriented Languages With Polymorphic Inline Caches. In *Proc. of ECOOP*, volume 512 of LNCS, pages 21–38. Springer, 1991. ISBN 3-540-54262-0.
- T. Kalibera, P. Maj, F. Morandat, and J. Vitek. A Fast Abstract Syntax Tree Interpreter for R. In *Proc. of VEE*, pages 89–102. ACM, 2014. ISBN 978-1-4503-2764-0.
- S. Marr, T. Pape, and W. De Meuter. Are we there yet? simple language implementation techniques for the 21st century. *IEEE Software*, 31(5):60–67, September 2014.
- T. A. Proebsting. Optimizing an ANSI C Interpreter with Superoperators. In *Proc. of POPL*, pages 322–332. ACM, 1995.
- A. Rigo and S. Pedroni. PyPy’s Approach to Virtual Machine Construction. In *Proc. of DLS*, pages 944–953. ACM, 2006.
- G. Sullivan. Dynamic Partial Evaluation. In *Programs as Data Objects*, volume 2053 of LNCS, pages 238–256. Springer, 2001.
- C. Wimmer and S. Brunthaler. ZipPy on Truffle: A Fast and Simple Implementation of Python. In *Proc. of OOPSLA Workshops, SPLASH ’13*, pages 17–18. ACM, 2013.
- T. Würthinger, A. Wöß, L. Stadler, G. Duboscq, D. Simon, and C. Wimmer. Self-Optimizing AST Interpreters. In *Proc. of DLS*, pages 73–82, 2012. ISBN 978-1-4503-1564-7.
- T. Würthinger, C. Wimmer, A. Wöß, L. Stadler, G. Duboscq, C. Humer, G. Richards, D. Simon, and M. Wolczko. One VM to Rule Them All. In *Proc. of Onward!*, pages 187–204. ACM, 2013. ISBN 978-1-4503-2472-4.
- A. Wöß, C. Wirth, D. Bonetta, C. Seaton, C. Humer, and H. Mössenböck. An Object Storage Model for the Truffle Language Implementation Framework. In *Proc. of PPPJ*, pages 133–144. ACM, 2014. ISBN 978-1-4503-2926-2.
- W. Zhang, P. Larsen, S. Brunthaler, and M. Franz. Accelerating Iterators in Optimizing AST Interpreters. In *Proc. of OOPSLA*, pages 727–743. ACM, 2014. ISBN 978-1-4503-2585-1.