

The Impact of Meta-Tracing on VM Design and Implementation

Carl Friedrich Bolz^a, Laurence Tratt^b

^a*Heinrich-Heine-Universität Düsseldorf, 40204 Düsseldorf, Germany.*

^b*King's College London, Strand, London, WC2R 2LS, United Kingdom.*

Abstract

Most modern languages are implemented using Virtual Machines (VMs). While the best VMs use Just-In-Time (JIT) compilers to achieve good performance, JITs are costly to implement, and few VMs therefore come with one. The RPython language allows tracing JIT VMs to be automatically created from an interpreter, changing the economics of VM implementation. In this paper, we explain, through two concrete VMs, how meta-tracing RPython VMs can be designed and optimised, and, experimentally, the performance levels one might reasonably expect from them.

Keywords:

Virtual machines, meta-tracing, programming languages.

1. Introduction

Every programming language that makes it beyond a paper design needs a corresponding implementation. Traditionally, most languages were compiled to machine code (via assembler or, less commonly, C). For languages with high static content (i.e. those with minimal runtimes) it can lead to highly efficient implementations, though it requires significant manpower to do so. However, languages with high dynamic content (i.e. those with complex runtimes, which includes most dynamically typed languages [1]) lack the static information needed to produce efficient code. In short, the traditional

Email addresses: cfbolz@gmx.de (Carl Friedrich Bolz), laurie@tratt.net (Laurence Tratt)

approach is generally either too costly, or leads to slow execution. For many modern languages, it is both.

It is now common for such languages to split apart the concept of compiler and ‘execution engine’, with the latter implemented as a Virtual Machine (VM). A simple compiler is all that is needed in such an approach, as the VM specifies most of the language’s behaviour. Implemented as a naive interpreter, a VM will tend to have poor performance. Instead, a VM can use information about the way the program executes to optimise it in ways that are statically impossible. Self first showed how a carefully crafted VM could substantially improve performance for a highly dynamic language [2] and its influence has led to most new languages being implemented using VMs. For example, programs running on Java’s HotSpot VM [3] (one of several Java VMs, henceforth referred to generically as ‘JVMs’), which is derived in part from the Self VM, can often match C’s performance.

However, a VM reflects the language, or group of languages, it was designed for. If a language fits within an existing VM’s mould, that VM will probably be an excellent target; if not, the *semantic mismatch* between the two leads to poor performance for user programs. For example, despite HotSpot’s excellent performance, Jython (Python for JVMs) almost never exceeds CPython (Python running with a custom VM written in C) in performance, and is generally slower, because features that make use of Python’s run-time customisability have no efficient implementation on a JVM. Similar problems have been observed when trying to implement Scheme on the JVM [4].¹ One attempt to alleviate this is the JVM’s `invokedynamic` bytecode [6], but special-casing never ends: if a runtime needs continuations or tail-calls, for example, the JVM has no natural way of expressing them.

Thus, languages which do not fit an existing VM mould must look to a custom VM to achieve good performance levels. However, this is not easily achieved: implementing a performant VM is justifiably seen as requiring highly specialized skills and significant manpower. In particular, the highest performing VMs rely on Just-In-Time compilers (henceforth referred to as JITs), which, at run-time, take part of a program and convert it to optimised

¹An early attempt at making a generic environment for implementing dynamic languages was done in the context of the Self project [5] where Java and Smalltalk bytecode was compiled to Self bytecode. Due to the Self VM’s excellent performance, the resulting translations performed similarly to the Smalltalk and Java VMs of the time. However, similar concerns about the wider applicability of the approach remain.

machine code. Few VM teams have the resources to create a JIT, particularly for complex languages. Unfortunately, therefore, most VMs' performance is substantially below that of HotSpot or .NETs CLR.

In this paper, we consider the RPython language, which allows automatic creation of JITing VMs from implementations of traditional interpreters through *meta-tracing*. This allows high-performance custom VMs to be created with reasonable resources. Many of the low-level details of RPython have been described elsewhere (see e.g. [7]). In this paper, we look at how a high-level VM is implemented in RPython, and the trade-offs involved, using two different VMs: PyPy [8] and Converge [9]. The two VMs have very different aims: PyPy is a drop-in replacement for the standard Python VM; Converge is a VM for a research language. Consequently, the two have had different levels of effort put into them (PyPy about 60 man months, not including the time spend to develop RPython; Converge about 3). Our aim in this paper is to explain, using Converge and PyPy as concrete examples, how RPython VMs can be designed and optimised, and what performance levels one might reasonably expect from them. This paper is the first to consider: specific RPython VM designs; the general lessons one can learn from them; and the effects of different man-power levels on such VMs.

Though the VMs we look at in this paper are for two relatively similar languages, it should be noted that RPython VMs have wider applicability. Pyrolog [10], an RPython VM for Prolog, is one example of RPython being used to implement a language that is not Python-like. Most of the techniques we discuss in this paper can be applied to a wide variety of languages; a few are specific to a given language or family of languages.

We start looking at the languages that each VM implements (Section 2), before introducing the RPython language itself (Section 3). RPython produces meta-tracing JIT, so we give an introduction to this area for the many readers who are likely to be unfamiliar with it (Section 4). We then detail the high-level design of the VMs themselves (Section 5), explaining them as if they were traditional interpreters. From that follows the technical heart of the paper, explaining the optimisation techniques the VMs use to generate performant RPython JITs (Section 6), as well as the general lessons embodied in the specific techniques. To show where RPython VMs with the optimisation techniques fit into the performance landscape, we present a performance comparison of various open-source VMs alongside the RPython VMs (Section 7). Finally, we look at the issues and limitations of tracing and meta-tracing JITs (Section 8).

The experimental suite for this paper is almost fully automated to enable repeatability: it automatically downloads and builds VMs, and then runs the experiments. We encourage interested readers to download and run it from http://tratt.net/laurie/research/publications/files/metatracing_vms/

2. Python and Converge

Python and Converge are two seemingly similar languages. Both are dynamically typed, object orientated, have an indentation-based syntax, and a rich collection of built-in datatypes. Indeed, much of Converge was explicitly influenced by Python. In this paper, we assume a passing familiarity with Python (pointing out differences with Converge on an as-needs basis).

Two technical features in particular distinguish Converge from Python. First, Converge allows compile-time meta-programming whereby code can be executed and generated at compile-time [9] (for the purposes of this paper, this can be thought of as approximately equivalent to Lisp macros). Second, Converge’s expression evaluation system is based on Icon’s and can perform limited backtracking. This second feature is interesting from the implementation point of view because it imposes a noticeable performance penalty even on programs which make little or no use of it [11].

Perhaps more importantly, at a ‘social’ level, the two VMs have different motivations. Python is a real-world language, used by hundreds of thousands of developers worldwide for a huge number of tasks, and for which many external libraries are available. PyPy’s goal is first to be fully compatible (warts and all) with CPython, and then to be faster. Converge on the other hand is a research language, intended to explore research on Domain Specific Languages (DSLs) and compile-time meta-programming. The needs of each language’s VMs reflect this: PyPy strives to be as fast as possible; Converge strives to be ‘fast enough’.

3. RPython

The basic facts about RPython are that it is a strict subset of Python whose programs are translated to C. Every RPython program is a valid Python program and can also be run using a normal Python interpreter. However, RPython is suitably restricted to allow meaningful static analysis. Most obviously, static types (with a type system roughly comparable to

Java’s) are inferred and enforced. In addition, extra analysis is performed e.g. to assure that list indices are not negative. Users can influence the analysis with `assert` statements, but otherwise it is fully automatic. Unlike seemingly similar languages (e.g. Slang [12] or PreScheme [13]), RPython is more than just a thin layer over C: it is, for example, fully garbage collected and has several high-level datatypes.

In addition to outputting optimised C code, RPython automatically creates a second representation of the user’s program. Assuming RPython has been used to write an interpreter for language L , one gets not only an optimised version of that interpreter, but also an optimising tracing JIT compiler for under 10 additional lines of code [14]. In other words, when a program written in L executes on an appropriately written RPython VM, hot loops (i.e. those which are executed frequently) are automatically turned into machine code and executed directly. As we shall see later, language implementers can influence the particular JIT that is created, using their knowledge of language semantics to allow further optimisations to occur.

RPython is able to automatically create JITs because of the particular nature of interpreters. An interpreter, whether it be operating on bytecode or ASTs, is a large loop: ‘load the next instruction, perform the associated actions, go back to the beginning of the loop.’ In order to switch from interpretation to JITing, RPython needs to know when a hot loop has been encountered, in order to generate machine code for that loop and to use it for subsequent executions. In essence, one need only add two annotations in the form of function calls to an RPython program to add a JIT. The first annotation informs RPython that a loop in the user program at position pc has been encountered, so it may wish to start generating machine code if that loop has been encountered often enough. The second annotation informs RPython that execution of the program at position pc is about to begin: it provides a safe-point for switching from the execution of machine code back to the interpreter.

4. Tracing JITs

Traditional JITs are *method JITs*: when a particular method is identified as being ‘hot’, it is translated into machine code (leaving most of its control structures intact).

In contrast, the JITs that RPython creates are *tracing JITs*. Tracing JITs came to prominence in the Dynamo project [15] as well as Franz and Gal’s

User program	Trace when x is set to 6	Optimised trace
if x < 0:	guard_type(x, int)	guard_type(x, int)
x = x + 1	guard_not_less_than(x, 0)	guard_not_less_than(x, 0)
else:	guard_type(x, int)	x = int_add(x, 5)
x = x + 2	x = int_add(x, 2)	
x = x + 3	guard_type(x, int)	
	x = int_add(x, 3)	

Figure 1: An example of a user program and resulting traces.

work [16]. The basic idea behind tracing JITs is to identify hot loops, record the bytecodes taken during a specific execution of it (‘the trace’), optimise the trace, and then convert that into machine code. Traces intentionally linearise control structures, naturally inlining functions. Wherever a specific branch was taken, a *guard* (roughly speaking, a ‘check’) is inserted into the trace; if, during execution of the machine code version, a guard fails, execution returns to the interpreter. Traces are hoped to be records of commonly taken paths through a program; when that assumption holds true, the result is extremely fast execution.

Figure 1 shows a high-level example of a program and its trace. The left-hand column shows a user program written in a Python-like language. When it is detected to be in a hot loop, the next time the code is executed, a trace is recorded. The middle column shows the trace recorded when x is set to 6. Note that the specific value of x is not recorded in the trace; indeed the trace would have been identical for any value of x greater than or equal to 0 (since the ‘else’ branch of the `if` would be taken for all such values); but the trace would be different if x was less than 0 (as the ‘then’ branch would be taken) or if x was not an integer. Once the trace has been recorded, the trace optimiser then attempts to reduce it in size, so that the resulting machine code executes as fast as possible. In this case, two type checks which are trivially true can be removed, and the two constant integer additions can be constant-folded. The resulting optimised trace is shown in the right-hand column.

4.1. Meta-tracing

Whereas tracing JITs are normally separate components from interpreters, RPython is a meta-tracing system [14]. The RPython translator in fact outputs two interpreters: the *language interpreter* is the (conceptually) simple

translation of the RPython interpreter into C; the *tracing interpreter* is a second representation of the interpreter which can be run to create traces. When a hot loop in a user program is detected, a marker is left such that the next time the loop is about to run, the VM will use the tracing interpreter instead of the language interpreter. When the loop is next encountered, a complete execution of the loop is performed and each low-level action taken by the tracing interpreter is recorded. After the loop has finished, the trace is then analysed, optimised, and converted into machine code. All subsequent executions of the loop will then call the machine code version. RPython automatically inserts guards into the machine code to detect divergence from the machine code version's capabilities. If a guard fails at any point, execution falls back to the tracing interpreter for the rest of that bytecode, and then back to the language interpreter.²

The fundamental difference between meta-tracing and non-meta-tracing JITs is that the latter JIT must be manually written. By tracing the actions the interpreter itself takes, a meta-tracing JIT can automatically create a JIT from the interpreter. As we shall see, the way an RPython interpreter is written affects the performance of the resulting JIT. To obtain the highest possible performance, the interpreter often needs to be subtly rewritten in specific places to aid the resulting JIT. When this is done intelligently, the raw traces created by an RPython JIT will often be reduced by 90% by RPython's trace optimiser [17].

The only other meta-tracing system we are aware of is SPUR [18], a tracing JIT for CIL bytecode, which can be used as a meta-tracer for languages implemented in C#. The sole paper on SPUR uses meta-tracing to implement a JavaScript VM (we are not aware of it being applied to other languages). The JavaScript interpreter is carefully structured so that the traces that SPUR produces for common object operations are efficient, yielding excellent performance results. This is similar in intent to RPython, though SPUR works at the C# bytecode level, and has fewer ways to annotate the interpreter to produce efficient traces.

5. PyPy and Converge VM overview

In this section we give an overview of the structure of both VMs.

²The reason that the tracing interpreter is only run sparingly is that it is extremely slow in comparison to the language interpreter.

5.1. Bytecode structure

Both PyPy and Converage use a bytecode based interpreter together with a compiler that translates programs into the respective bytecode set. The bytecode sets are similar in intent, being stack-based and deferring type specialization until run-time. Python's bytecode set contains more instructions to optimise specific common operations (e.g. list accesses) than Converage, while the latter has several additional instructions related to backtracking.

5.2. Compilation

Both VMs store programs as bytecode for execution by the eventual interpreter. Although both PyPy and Converage use traditional compilation, the implementations differ. PyPy's compiler is written in RPython and is integrated into the VM for fast startup times; most users will never be aware that separate compilation is performed on their behalf. Converage's compiler is written in Converage, and is a separate program; while it can be invoked manually, if passed a source file, the VM transparently calls the compiler. Both systems attempt to transparently cache bytecode output on disk to lower compiler costs, automatically recompiling any source files which are detected to be newer than their cached equivalents (this means that the first run of a Python or Converage program can be significantly slower than subsequent runs).

Because the Converage VM is used to compile new versions of the Converage compiler, the latter has to obey an important restriction: neither the compiler nor any libraries it uses can perform compile-time meta-programming. If the compiler were to do so, it would be impossible to migrate Converage's bytecode format, as the running compiler would then emit bytecode in the new VM format and attempt to execute it, all while still running on the old VM. In practice, this restriction is not particularly onerous, although it requires a freshly unpacked Converage system to be compiled in a specific order: first a minimal version of the standard library (enough for the compiler); then the compiler itself; then the full library (which, at this point, may include compile-time meta-programming).

5.3. Interpreter structure

Both PyPy and Converage split their interpreters into three major parts: the bytecode interpreter; the built-in datatypes; and the built-in libraries. The bytecode interpreters are responsible for dispatching and implementing the bytecode set and are constructed in a direct, simple fashion. The built-in

datatypes realise basic concepts such as objects, classes, lists, and dictionaries. As well as being used extensively throughout the VM, several of these datatypes require careful bootstrapping during VM initialization. Built-in libraries are provided either for performance reasons or to allow integration with low-level C libraries.

In Converge, the split between the bytecode interpreter and built-in datatypes is relatively informal, as befits a simple VM. In PyPy, in contrast, the split is very clearly defined to ensure that, despite the large size of the Python language specification, the components are manageable. The bytecode interpreter treats all Python objects that it handles as black boxes; operations on them are handled by a separate component called the *object space*. The only way for the bytecode interpreter to gain actual knowledge about an object is to ask the object space for the object's truth-value (i.e. whether the object is equivalent to True or False, information necessary for `if` statements). The object space, on the other hand, only knows about datatypes, not about executing Python code, for which it refers back to the interpreter.

High-level languages such as Python typically have a Foreign Function Interface (FFI) to interface to external C libraries. Because of the mismatch between the high-level language and C, FFIs are often clumsy to use. RPython's more static nature and lower-level types make it a better fit: consequently, PyPy and Converge mostly interface to C libraries in RPython.

Libraries which do not need to interface to external C libraries are more interesting in an RPython VM. Traditional VMs such as CPython implement as much functionality in C as is practical, often migrating libraries from the native language to C over time. The speed advantages of doing so are often huge, and such language communities develop careful conventions about what calculations should be done via library calls to take advantage of this. An important goal of RPython VMs is to significantly reduce the need to write and use C modules for performance reasons.

6. Optimising an RPython VM

Optimising an RPython VM means concentrating on the two execution modes: optimising the interpreter for faster interpretation speed; and rewriting the interpreter to produce traces which can be better optimised by the JIT. The former is largely similar to the challenges faced by other interpreters, so we dwell little on it; the latter is more unique to RPython VMs

and what we concentrate on in this section.

From the perspective of an RPython VM author, many standard optimisations ‘fall out of the hat’. Built-in datatypes such as integers, floats, and (to an extent) strings are naturally optimised by RPython’s allocation removal techniques [17].

What an RPython VM author needs to concentrate on are the commonly used building blocks that are specific to the language being implemented. In the case of Converge and PyPy, the three common pinch points are instances (objects), classes, and modules.³ As highly dynamic languages, Converge and Python programs can change and inspect run-time behaviour in arbitrary ways. However, most programs restrict such changes to small portions. Both RPython VMs therefore aim to make the common case of non-reflective access as fast as possible. Conversely, when a program uses the language’s more dynamic features (introspection, self-modification, intercession [19]), execution falls back to slower, more general code.

In this section we give an overview of how language-specific building blocks can be optimised; many of the techniques described will be applicable to the different building blocks found in other languages.⁴ In general, the Converge VM implements the ‘easy win’ optimisations, while PyPy optimises a much wider class of programs. Both experiences are useful: Converge shows how significant optimisations are possible with little effort, while PyPy shows how RPython VMs can optimise seemingly resistant programs.

6.1. General RPython JIT optimisation techniques

The techniques described in this section are more finely-tuned variants of the techniques described in [21]. The general aim is to produce small traces which can then be further shrunk by RPython’s trace optimiser. The overall strategy is to expose, by rewriting the interpreter, the parts which can be made constant in traces based on that code; these parts can then be optimised away, leaving only simple guards in their place. The tactics used to achieve this involve either using RPython-level annotations (i.e. promoting values and eliding functions) or rewriting the interpreter to use more trace-friendly code (e.g. moving from arbitrarily sized arrays to fixed-size lists). We now give a brief explanation of each.

³Informally, in both Converge and Python, a ‘library’ is a collection of modules.

⁴For PyPy container types (e.g. lists and hash maps) optimisations, see [20].

6.1.1. Promoting values

In the context of a specific trace, it is often reasonable to assume that certain pieces of information are constant. The trace optimiser can be informed of this likelihood by *promoting* a value. For the small cost of inserting a guard at the first point of use, all subsequent calculations based on that constant can be constant-folded or removed from the trace. Note that constants are not known at compile-time: they are run-time values that are only constant for one particular trace. An important example of this is the concrete type of an object. Even in dynamically typed languages, most variables are only ever assigned values from a small set of types. Promoting the type of an object allows calculations to be specialized on that type. Because there is a very high likelihood that only a single type will be used at a given program point, the corresponding guard will fail only rarely.

6.1.2. Elidable functions

Similarly to promoting a value, functions can be annotated as being *elidable*. This is similar, though not identical, to the concept of pure functions. In short, an elidable function must guarantee that, given the same inputs, it always returns the same outputs. In contrast to pure functions, elidable functions may also have side effects (e.g. caching), provided that the same inputs always result in the same outputs being returned. When a call to such a function is encountered in a trace, its body thus need not be executed when the input values match those previously encountered.

6.1.3. Using trace optimiser friendly code

Often seemingly similar techniques can yield surprisingly different results in the context of the trace optimiser: one might frustrate the optimiser; another may allow it do its job well. As a concrete example, we look at the most common collections datatype: lists.

Arbitrarily resizable lists cause the trace optimiser something of a headache. Every `append` (or similar) operation requires a check to see if the list has enough space left; if not, it must be resized, and possibly moved elsewhere in memory. Because of this, the trace optimiser can not definitively prove useful properties of the list over the lifetime of a trace, and therefore can not optimise much. Whenever possible, therefore, arbitrarily resized lists should be avoided. We now give two indicative solutions to avoiding the use of such lists.

Fixed sized arrays are the most obvious solution. These are much more amenable to trace optimisation as they expose constant information – in this case, the size of a list – to the trace optimiser. Indeed, the single biggest improvement in performance in the Converge VM was moving from a global stack (as a resizable list) to a per-function frame stack (as a fixed size array). This necessitated modifying the compiler to statically calculate the maximum stack space a function requires at run-time (roughly one day’s work), and creating a fixed-size list of that size in each function frame.

Arbitrarily sized lists which are not randomly accessed need not be stored contiguously at all, instead being accessed as a linked list. A simple example of this is function frames. Since one needs to be able to access all of these to print out backtraces, an old version of the Converge VM, stored these contiguously in an arbitrarily resizable list. Clearly we can not replace such a list with a fixed size array: we have no idea in advance how deep functions will recurse. However, we do know that we rarely need to access anything other than the current frame’s parent (to know where to return at the end of a function call). Therefore, having each frame store a pointer to its parent frame is a simple solution. It gives quick access to the parent frame, and, via pointer traversal, all the grandparent frames when a backtrace is needed. However, in essence the explicit list has disappeared entirely, and the trace optimisers life becomes much easier.

6.2. *Optimising Instances*

Both Converge and Python allow users to define classes and create instances (i.e. new objects) from them. Most programs create many such instances. Optimising the common cases is thus extremely important. We also use instances as an example of how rarer, but more complex, language semantics can be handled without affecting the common case.

The basic semantics are similar in Converge and Python. Every instance records the class it instantiates; both languages allow this to be changed at run-time (in Python by writing to the `__class__` slot, in Converge to `instance_of`). Instances can also store an arbitrary number of slots (key/value pairs), which can vary on a per-instance basis (i.e. unlike many other OO languages, a class does not precisely define the ‘shape’ of its instances). In essence, instances behave like dictionaries mapping slot names (as strings) to values, while classes define the shared behaviour between instances. In this sense, both languages behave more like prototype-based languages such as Self than class-based languages such as Smalltalk.

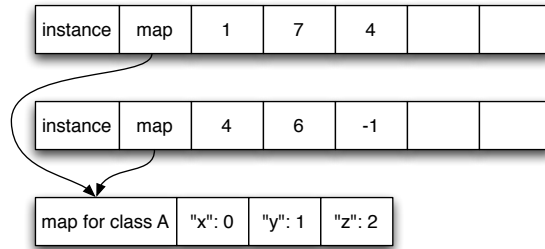


Figure 2: Two instances of class A sharing the same map

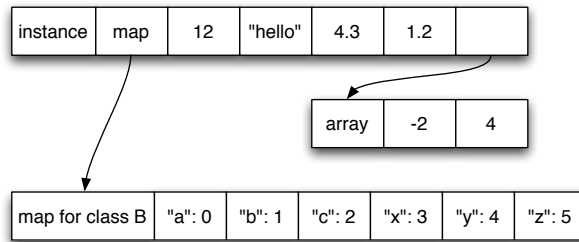


Figure 3: An instance of class B with six slots

Both Converge and PyPy optimise the common case of directly accessing slots in instances using *maps* (a concept originally from Self [22]; the technical details of PyPy’s approach to maps are explained in [21]). Although instances of the same type can vary substantially, in practice the ‘shape’ of an instance is highly correlated with its class. Since the two rarely vary independently, PyPy stores the references to the class of an instance in its map, not directly in the object, saving a promotion in the process. Since promotions turn into guards in a trace, this produces smaller traces.⁵ It also has the benefit of making objects one word smaller.

Although performance is PyPy’s most obvious goal, it also attempts to save memory when that is not in direct conflict with performance. One example of this is PyPy’s compact representation of instances. An informal study of real systems showed that most objects have 5 or fewer slots. PyPy therefore preallocates space for 5 slots, freeing it from the need to allocate an arbitrarily sized list to store slots in most cases (which, when all of its parts

⁵Of course, many trace operations might be later optimised away; in this case, the guard resulting from the second promote would be unlikely to be so optimised.

are taken into account, needs around 40% more memory to store 5 slots). Only when more than 5 slots are assigned to an instance is an arbitrarily sized list created and referenced from the object.

Figure 2 shows PyPy’s layout scheme for two instances of class A, each instance using the same additional slot names. Since the instances have only three slots, the content of the slots can be stored in the free slots. Figure 3 shows an instance with six slots. Two of the fields have to be stored in an extra array allocated for that use. Note how the last field of the instance is used for the indirection.

6.2.1. Python’s additional instance semantics

Python’s instance model has a number of complexities over Converge’s, which PyPy fully supports. These complexities are interesting because they show how interpreters can gradually allow performance to tail off as rarer, more dynamic, features are used.

The complexities relate to an implementation decision from CPython: every instance has a reference, via the `__dict__` slot, to a dictionary that stores all the instance’s slots. This dictionary can be replaced by writing to the `__dict__` slot, changing all the instance’s slots. This implementation decision is costly in terms of memory, as dictionaries are not small data structures, and seems to defeat many reasonable optimisations.

One solution would be to use maps for normal accesses, but switch to a plain dictionary as soon as the `__dict__` slot is accessed. Doing so would mean that any reflective access of the dictionary would slow down all subsequent uses of that instance. Since the dictionary is mostly used for reading and writing slots, this would slow down many real programs. Therefore, in PyPy, requesting an instance’s dictionary returns a fake dictionary. This is indistinguishable from a real dictionary, and transparently redirects all reads and write to keys and values to the underlying instance.⁶ In other words, performance for normal accesses remains as fast as the standard case.

Figure 4 shows an instance, its map, and the fake dictionary that redirects all accesses back to the instance. Note that the instance needs to keep a

⁶A more complete solution for this sort of reflective access would be to use mirrors [19]. However, this would require changing the semantics of the Python language. In some senses, the `__dict__` attribute can already be seen as a mirror on the attributes of an instances. Indeed, it gives additional guarantees over mirrors, guaranteeing that the identity of `__dict__` is the same on all accesses.

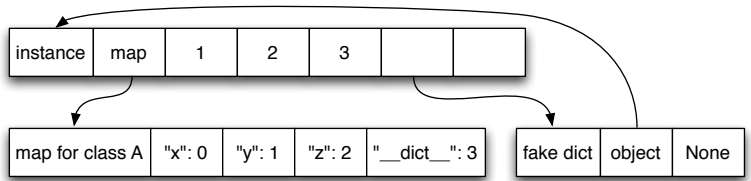


Figure 4: An instance implemented with a map, and its dictionary

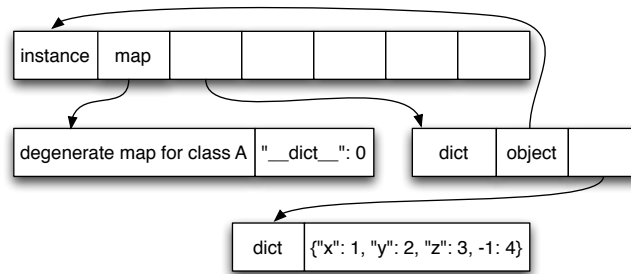


Figure 5: An instance that has its slots stored in a dictionary

reference to the dictionary once it has been requested in order to ensure that the expected object identity invariants are maintained.

However, when the programmer uses more of Python’s dynamic features – in particular, writing a new dictionary to the `__dict__` slot – even this tactic is no longer viable. In such cases, PyPy stops using maps for the instance and stores its instances in a real dictionary (as shown in Figure 5). Fortunately such uses are rare, so few programs suffer the consequent slowdowns.

6.3. Optimising Classes

Both Python and Converage instances store only the information which varies from the class they instantiated. Typically this means that instances store dynamic information (ints, strings, user classes etc.) while classes store static information (typically functions). Accessing fields in classes is thus as common an operation as accessing slots in instances. Both PyPy and Converage aim to make non-reflective method lookup as fast as possible.

Looking up a method in a class necessitates, conceptually, traversing its inheritance hierarchy (note that both Converage and Python support multiple inheritance; Python uses the C3 algorithm [23], which means potentially looking at all its base classes during every method lookup). Since both languages allow classes to change dynamically, method lookup is a seemingly expensive operation.

The technique both languages use is to *version* classes. Every change to a class (e.g. adding or editing a field) changes its version. For any given version of a class, all of its fields are thus constant, and accesses to that class can be promoted (based on both the class *and* the version) and elided away. Because of inheritance, classes can not be versioned in isolation: for example, if a field is added to a class, then instances of its subclasses should gain that field too. Thus, as well as changing the version of a class when it is edited, we must change the versions of each of its subclasses. Since storing a normal reference to subclasses would prevent the latter ever being garbage collected, both Converse and PyPy classes store weak references (i.e. references that do not keep their target object ‘alive’) to their subclasses.

This technique makes looking up a field in a class extremely quick (comparable in speed to C++ method calls) for the common case. The JIT optimises field lookups to a single guard which need only check that one class’s version; if the check succeeds, the correct result is already known and inserted. Since versions can change an unbounded number of times, the seemingly obvious technique of using a monotonically incrementing integer for the version is dangerous: the integer could then overflow and two versions that were intended to be different could appear to be the same version, leading to unexpected behaviour. However, we only need compare whether one version is different than another, not whether one version is newer or older than another. Both PyPy and Converse therefore instantiate blank objects of an arbitrary class to stand in for versions: the RPython memory system implicitly guarantees that two different objects will compare differently, providing exactly the guarantees needed, without any possibility of integer overflow.

However, as presented above, performance would suffer for the rarer case where a class’s fields change frequently such as when a class stores a monotonically increasing counter which it assigns to every instance. Every instantiation would change the class and its subclass’s versions; worse, trace guards would be invalidated and traces begun anew.

PyPy therefore adds one technique to the above (which Converse does not currently do). When a class field is given a different value for the first time, an extra level of indirection is introduced: the class no longer stores the field’s value directly, but stores a reference to a small intermediate object (a class cell) that contains the value. When that particular field is changed subsequently, only the content of that object is changed, not the class as a whole: the class’s version therefore need not be changed. After the first time, writing to such a field causes relatively little slowdown, while reading from it

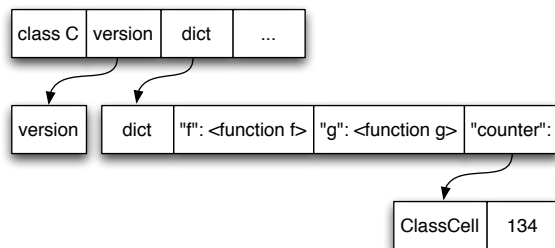


Figure 6: Class C with two methods and a counter

needs an extra memory read (including when accessed via subclasses). While slightly less efficient, this balances fast general performance with reasonable performance in the rarer case.

Figure 6 shows a class with two methods `f` and `g` and a `counter` field. The counter is stored via an indirection to a `ClassCell`, so that changing it does not update the version of the class. Reading the `counter` slot requires an extra pointer dereference.

This approach is similar to Smalltalk’s mechanism for handling global variables [24, p. 599]. A global variable⁷ is a reference to an *association* object, which corresponds to PyPy’s class cells. To read it, the value of the association is read. To write to the global variable, the value of the association object is set. In Smalltalk, this indirection is always used, not just for commonly changed variables. More generally, class versions can be related to the invalidation of method caches when a new method is compiled in some Smalltalk systems [25].

6.4. Optimising Modules

Modules are conceptually similar to classes, with both providing namespaces for storing functions/methods and values. However, modules are simpler in that there can be no inheritance between them. Modules in PyPy therefore use a similar versioning technique to classes. Converge uses a more naive scheme, to maintain simplicity in the compiler and VM. Top-level Converge module scopes are simply closures and, within a module, they can be assigned to as normal; synchronising their mutation within and without the module would be somewhat difficult. Furthermore, tracking the number of

⁷The same is true for class variables and pool variables.

assignments and adding indirection would be another complication. Converge modules thus use maps (which are promotable, and calculations on them easily elided) to map module lookups and assignments to an offset in a closure (which is a fixed size array). In practice, all reads and writes to Converge modules act like indirected accesses in PyPy. This gives reasonable (though not stellar) performance with little effort.

6.5. Discussion

With the optimisations described in this sections so far, instances, classes, and modules perform well in both the Python and the Converge VM. Instances are stored almost as compactly in memory as HotSpot, with equally efficient attribute access times, despite retaining sufficient information to implement highly dynamic languages. Classes are highly optimised for the common case (an inheritance hierarchy where methods in classes are not changed). In PyPy, module globals have most of their lookup overhead removed; in Converge, they are less efficient, but still adequately fast.

These optimisations exemplify how RPython VM authors need to consider which usage patterns are the most important (i.e. frequent) and therefore should be made as efficient as possible. They must then (re)arrange the interpreter and data structures so that, in conjunction with the trace optimiser, small traces with little code and few guards are produced. There is, of course, a tension between making common cases fast while not making less common cases unusably slow. VM authors need to understand their languages and intended use cases well. However, as often the case with performance issues, it is not realistic to do so purely intellectually: real programs must be analysed to determine which cases need to be focused on. Different benchmarks (synthetic or not) can change the perception of the most important areas substantially, and must be carefully chosen.

As this suggests, it is impossible to design a perfectly optimal interpreter up-front. Analysing traces from real programs often shows new opportunities for optimisation. Each pinch-point identified in the interpreter can be addressed either by adding hints for the JIT, or by rewriting the interpreter. This is often not a trivial task, particularly for more complex interpreters. It requires careful thought about the goals of the optimisation, the trade-offs involved (including to code readability), and how to reach these goals.

Thus, while a basic meta-tracing JIT comes ‘for free’, a fully optimised one is no small task. That said, nearly all optimisations are understandable at the level of the interpreter itself: one need never look within the JIT

compiler itself. The interpreter thus still expresses the language semantics correctly – albeit somewhat strangely when optimisations require changing its structure – and many optimisations improve the performance of the language interpreter as well as the resulting JIT. For example, maps are a memory optimisation even if only an interpreter is used, and version tags can be used for a method cache within a purely interpreted system. Such rewritings are akin to an extreme version of rewriting a C program knowing the sorts of idioms that are best optimised by the C compiler.

While promoting and eliding are direct features of the JIT, version tags are an idiom of use. This can understate their importance: they are a powerful way to constant-fold arbitrary functions on large data structures. The versions need to be updated carefully every time the result of a function on a structure can change. Therefore this technique is only applicable on data structures which change slowly or which (as PyPy’s approach to optimising classes in Section 6.3 showed) can be made to change slowly.

We believe that the manual rewriting of parts of the interpreter is a key part of the meta-tracing approach. Many optimisations rely on in-depth knowledge of the language the interpreter implements. The rewrites expose not only properties of the language semantics (which are already present in the interpreter) but also expectations about patterns of use (which are not). While an ‘optimally smart’ meta-tracing compiler might deduce some optimisations, many require a human’s understanding of the wider context.

7. Performance

In this section we aim to show where RPython VMs fit in the general VM performance landscape. We therefore compare the PyPy and Converge VMs to: C (to show how hand-written programs compare to VMs in general); hand-crafted high-performance VMs (HotSpot and LuaJIT); hand-crafted low-performance VMs (CPython, Lua, Ruby, and Converge1); and translations to another VM (JRuby and Jython). In order to do this, we need programs which are available in each language. Synthetic benchmarks are the only plausible candidates for cross-language comparison. With a reminder to readers of the inevitable limitations of synthetic benchmarks – they can easily be ‘gamed’ by language implementers and are often not representative of real workloads – we explain the systems under test, the methodology we use to measure performance, and the experimental results.

Language implementation	KLoC	VM Impl. Lang	Language version
C (GCC 4.6.3)			
HotSpot (1.7.0_09)	250	C++	Java (1.7)
Converge1 (git #68c795d2be)	11	C	Converge (1.2)
Converge2 (2.0)	4	RPython	Converge (1.2)
Lua (5.2.1)	14	C	Lua (5.2)
LuaJIT2 (2.0.0)	57	C	Lua (5.1)
CPython (2.7.3)	111	C	Python (2.7.3)
Jython (2.5.3)	63	Java	Python (2.5)
PyPy-nonopt (1.9*)	31	RPython	Python (2.7.2)
PyPy (1.9)	33	RPython	Python (2.7.2)
Ruby (1.9.3-p327)	102	C	Ruby (1.9.3)
JRuby (1.7.1)	115	Java	Ruby (1.9.3)

Table 1: The language implementations we compare.

7.1. Systems Under Test

The benchmarks we use are: Dhrystone [26], a venerable integer benchmark, and almost certainly the most widely ported cross-language benchmark; Fannkuch-redux, which counts permutations [27], from the Computer Language Benchmarks Game⁸; and Richards⁹, which models task dispatch in an operating system. Dhrystone is included for its ubiquity; Fannkuch-redux for its exercising of built-in datatypes; and Richards for its relative real-worldism. In the Appendix we present 8 other benchmarks.

Table 1 shows the language implementations we compare, with detailed version information to ensure repeatability. Converge1 (the old C VM) is the Converge 1.2 VM with the minimal number of functions (related to integers and strings) added to allow the benchmarks to run (we give the git hash to allow precise recreation of this version). Converge2 is the new RPython VM. PyPy-nonopt is a variant of PyPy with the major optimisations of Section 6 turned off (the flags used to obtain this can be found in our repeatable build system; we call this version 1.9* to differentiate it from standard PyPy), allowing us to explain the effect of those optimisations.

We give Lines of Code (LoC) rounded to the nearest 1000 LoC – excluding blank and / or comment lines – for the ‘core’ of each VM¹⁰ to give an indica-

⁸<http://shootout.alioth.debian.org/>

⁹<http://www.cl.cam.ac.uk/~mr10/Bench.html>

¹⁰HotSpot numbers from <http://openjdk.java.net/groups/hotspot/>.

tion of relative size. As is commonly the case with LoC, the precise number should be treated with caution: different languages and coding styles make precise comparisons impossible; and VMs vary considerably in the extent to which library functionality is included in the VM or as a normal user library. With some VMs, determining which files should be counted as part of the VM or not is a matter of considerable debate. Nevertheless the LoC count allows one to get a rough handle on the effort levels that have gone into each VM. Both of the RPython implementations are about one third the size of their C counterparts. The optimisations described in Sections 6.2, 6.3, and 6.4 add about 2K LoC to PyPy. Neither the PyPy nor the Converge2 numbers include the RPython infrastructure or the meta-tracing JIT compiler.

Our experimental system is almost fully automated to enable repeatability: it automatically downloads and builds the correct versions of the VMs (except for HotSpot), calculates the LoC for each VM, and then runs the experiments. The download link can be found on page 4.

7.2. Methodology

A fundamental problem when measuring JIT-based systems is whether to include warm-up time or not. JIT implementers often argue that warm-up times are irrelevant for long-running processes, and should be discounted. Others argue that many processes run for short time periods, and that warm-up times must be taken into account. We see merit in both arguments and therefore report two figures for each benchmark: *short*, where the benchmark has a low input size (e.g. 10 for Richards), and where warm-up times can play a significant part; and *long*, where a higher input size (e.g. 100 for Richards) tends to dominate warm-up times.

We ran all systems using the default options, with 3 exceptions. First, we used the `-O3` optimisation level for GCC. Second, we increased the memory available to the HotSpot-based VMs, as otherwise several of the benchmarks run out of memory. Third, we used the `-Xcompile.invokedynamic=true` option for JRuby to force the use of HotSpot's `invokedynamic` instruction, which is otherwise disabled on current versions of HotSpot.

We ran each version of the benchmark 30 times using `multitime`¹¹ to randomise the order of executions on an otherwise idle Intel Core i7-2600S 2.8GHz CPU with 16GB RAM, running i386 Linux 3.5.0 and GCC 4.6 as the

¹¹<http://tratt.net/laurie/src/multitime/>

	Dhrystone		Fannkuch Redux		Richards	
	50000	5000000	10	11	10	100
C	0.004 ± 0.002	0.179 ± 0.010	0.163 ± 0.006	1.992 ± 0.010	0.012 ± 0.006	0.079 ± 0.006
HotSpot	0.107 ± 0.006	0.240 ± 0.010	0.350 ± 0.008	3.448 ± 0.029	0.109 ± 0.010	0.169 ± 0.014
Converge1	2.053 ± 0.029	207.274 ± 3.048	-	-	9.931 ± 0.102	100.216 ± 1.356
Converge2	0.118 ± 0.004	1.914 ± 0.022	2.658 ± 0.041	33.484 ± 0.517	0.637 ± 0.006	2.850 ± 0.014
Lua	0.201 ± 0.008	19.417 ± 0.474	7.683 ± 0.321	100.536 ± 2.475	0.665 ± 0.024	6.574 ± 0.139
LuaJIT	0.014 ± 0.006	0.879 ± 0.016	0.339 ± 0.008	4.180 ± 0.010	0.085 ± 0.006	0.763 ± 0.010
CPython	0.368 ± 0.010	35.072 ± 0.537	9.167 ± 0.237	114.001 ± 2.189	1.585 ± 0.022	15.698 ± 0.227
Jython	1.820 ± 0.029	28.432 ± 0.466	7.776 ± 0.419	76.069 ± 4.753	2.820 ± 0.069	13.870 ± 0.345
PyPy-nonopt	0.127 ± 0.006	5.898 ± 0.071	1.402 ± 0.022	16.989 ± 0.220	0.515 ± 0.010	2.839 ± 0.016
PyPy	0.069 ± 0.008	1.085 ± 0.014	1.256 ± 0.024	15.239 ± 0.223	0.267 ± 0.006	0.544 ± 0.008
Ruby	0.312 ± 0.008	29.819 ± 0.257	13.152 ± 0.200	172.098 ± 2.168	0.793 ± 0.018	7.159 ± 0.061
JRuby	2.050 ± 0.039	10.576 ± 0.304	6.313 ± 0.127	61.934 ± 1.513	2.130 ± 0.025	3.640 ± 0.053

Table 2: Benchmark Results.

compiler (full details can be found on the experiment website). We report the average wall time and confidence intervals with 95% confidence levels.

7.3. Experimental results

Table 2 shows our experimental results. Several things are worthy of note.

The interpreter-only VMs (Converge1, CPython, Lua, and Ruby) show similar, typically linear, slowdowns as the benchmarks lengthen. Lua is faster than CPython and Ruby, which are both faster than Converge1.

Jython and JRuby are both compilers which create JVM bytecode which runs on HotSpot. JRuby is generally faster than the Ruby interpreter while Jython is generally slower than CPython. JRuby is faster than CPython, probably due to its extensive use of the `invokedynamic` bytecode (see page 2). Given that both run atop HotSpot – which, on its own, is nearly always the fastest VM by a considerable margin – it may seem surprising that both Jython and JRuby are outperformed by the much simpler Converge2. We believe this underlines the ‘semantic mismatch’ problem we outlined on page 2.

Of the other JITted VMs, LuaJIT clearly outperforms PyPy and Converge2. This is most noticeable on the fast benchmark runs, which show that RPython JITs warm-up rather slowly (see Section 8). Since Lua is a small language, LuaJIT has been carefully hand-crafted for performance. Python is a significantly larger language, with many more complex corner-cases. We suggest this explains why similar hand-crafted Python JITs in the past (e.g. Psyco) have not been able to speed up all of Python’s features.

RPython has allowed PyPy to reach feature parity with relative ease—the performance trade-off therefore seems a reasonable compromise.

Converge2 is dramatically faster than its non-JITted predecessor Converge1 (and, indeed, the interpreter-only VMs Lua and CPython), but can not compete with the more carefully tuned PyPy. Converge1 has a memory corruption bug which shows up in the Fannkuch-redux benchmark, but which had never been noticed before. Converge2 has no such problems, relying on RPython’s automatic memory management.

Comparing PyPy–nonopt with PyPy shows that the optimisations of Section 6 benefit many programs, often substantially. Dhystone benefits from the optimisations because it uses many global variables; Richards because it’s written in a strongly object-orientated style. However, not all programs benefit. Since Fannkuch Redux mostly manipulates list and performs arithmetic, it has few code-paths which benefit from Section 6’s optimisations.

8. Issues

As Section 7 shows, RPython VMs typically exceed hand-written interpreters in performance, even when less effort has been put into them. However, it would be foolish to pretend that RPython VMs are without issues. In this section, we enumerate the issues specific to RPython, and those common to all tracing approaches, as well as giving suggestions for possible solutions.

The most obvious problem with RPython VMs is the time it takes to warm-up the JIT (i.e. for all the ‘hot spots’ in the code to have been traced and converted to machine code). Though all JITs suffer from this problem, the warm-up penalty in RPython VMs is larger because the JIT is language independent. Rather than having a custom tracer, the tracing interpreter (see Page 6) is used to generate tracers. The tracing interpreter is, in effect, itself interpreted to produce traces, causing a double interpretation overhead. This is then compounded by the fact that meta-tracing inevitably creates large traces, which are expensive to produce, analyse, and optimise.

Since tracing is an expensive activity, loops in an RPython VM must be executed many more times than a traditional JIT before tracing is started. Long-running processes take longer to ‘warm-up’ than might be expected, and short-running processes often derive little or no benefit from the JIT.

Solving this problem would involve replacing the automatically created tracing interpreter with an equivalent component that is able to produce

traces more efficiently. A plausible approach would be to specialize the tracing component to the language being implemented, removing the double interpretation overhead. The RPython project has previously experimented with this approach, but no part of the implementation remains.

The RPython language itself is not without issues. Some – such as a relative lack of documentation – are likely to be solved with time, and are not important enough to mention in detail here. A more fundamental problem is the ‘Python’ part of RPython. First, RPython is statically typed, but types can not be directly expressed in the language: they are instead inferred, with corresponding problems when inference goes awry. Second, RPython uses Python as its compile-time meta-programming language (roughly speaking, the language it uses to generate code at compile-time). The RPython translator loads in a normal Python file and executes it for as long as it chooses. Once that has finished, the translator expects to be given a reference to the VM’s entry point, whence translation occurs. Everything referencable from the entry point must be ‘RPython enough’ to be translatable; things not reachable are ignored (and may use arbitrary Python features). Compile-time meta-programming is vital for software that needs to be customisable and portable. However, it means that most VM files are in fact mixed Python and RPython programs. This mixing and matching of two similar, but distinct languages, is often confusing. Furthermore, it makes it difficult to translate RPython VMs in a modular fashion. Currently translation is ‘whole program’, and must be done for every single change. Large systems such as PyPy can take 45-60 minutes when a JIT is generated.

We suspect that future RPython-esque systems will choose a language with explicit static typing, and a clearly delineated compile-time meta-programming phase. For the former, a Java-esque language is likely to be sufficient; for the latter, a Converge-esque approach may yield good results.

A problem common to all tracing JITs (including those of RPython VMs) is that they give uneven performance improvements. Some programs run faster than method-based JITs while some run substantially slower. Programs in the latter category are invariably those which change the control flow paths they follow frequently. This causes guards to fail more often, and extra traces to be triggered. Compilers are a classic example of such programs (AST walkers appear, to a tracing JIT, to take different paths almost at random). The tracing JIT’s overhead can outweigh its benefits unless such programs run for a long time and all the common paths are traced.

From the point of view of RPython VMs, this problem could only be

solved by moving beyond the tracing paradigm. However, it is not clear how this might be done. First, tracing is a pragmatic way of getting reasonably good results for a wide variety of languages: other approaches are much harder to control and ‘tune’. Second, as RPython shows, tracing is particularly amenable to ‘meta’ approaches: it would be harder to automatically create a method-based JIT in this way, for example.

9. Conclusions

By looking at two different RPython VMs, we hope the reader has gained an understanding of the power of meta-tracing and its performance characteristics. We also detailed general optimisation techniques for meta-tracing VMs (as embodied in the PyPy and Converge VMs) that are likely, directly or indirectly, to aid future authors of meta-tracing VMs. We believe that further research into this area is likely to continue to narrow the performance gap with hand-crafted JITs. For those prepared to pay the high manpower costs, hand-crafted JITs will always retain a performance edge; however, as this paper has demonstrated, language implementations can now perform at reasonable performance levels with surprisingly little effort.

Acknowledgements: We thank Lukas Diekmann, Samuele Pedroni, David Schneider, Naveneetha Vasudevan, and the anonymous reviewers for insightful comments on drafts of the paper. Fabio Mascarenhas, Takafumi Nose, and Martin Richards kindly placed the Lua Richards, Ruby Dhrystone, and Richards benchmarks respectively into the public domain. We thank the PyPy and RPython community for their continuous support and work: Armin Rigo, Maciej Fijałkowski, Alex Gaynor, and countless others. Any remaining errors and infelicities are our own.

References

- [1] L. Tratt, Dynamically typed languages, *Advances in Computers* 77 (2009) 149–184.
- [2] C. Chambers, D. Ungar, Customization: optimizing compiler technology for SELF, a dynamically-typed object-oriented programming language, in: *Proc. PLDI*, ACM, 1989.
- [3] M. Paleczny, C. Vick, C. Click, The Java HotSpot server compiler, in: *Proc. JVM Research and Technology Symposium*, USENIX, 2001.

- [4] B. P. Serpette, M. Serrano, Compiling Scheme to JVM bytecode: a performance study, *SIGPLAN Not.* 37 (2002) 259–270.
- [5] M. Wolczko, O. Agesen, D. Ungar, Towards a universal implementation substrate for Object-Oriented languages, *Proc. Workshop on Simplicity, Performance, and Portability in Virtual Machine Design*, 1999.
- [6] J. R. Rose, Bytecodes meet combinators: invokedynamic on the JVM, in: *Proc. VMIL*, ACM, 2009.
- [7] D. Ancona, M. Ancona, A. Cuni, N. D. Matsakis, RPython: a step towards reconciling dynamically and statically typed OO languages, in: *DLS*, ACM, 2007.
- [8] A. Rigo, S. Pedroni, PyPy’s approach to virtual machine construction, in: *DLS*, ACM, Portland, Oregon, USA, 2006.
- [9] L. Tratt, Compile-time meta-programming in a dynamically typed OO language, in: *Proc. DLS*, ACM, 2005, pp. 49–64.
- [10] C. F. Bolz, M. Leuschel, D. Schneider, Towards a jitting VM for Prolog execution, in: *PPDP*, ACM, Hagenberg, Austria, 2010.
- [11] L. Tratt, Experiences with an Icon-like expression evaluation system, in: *Proc. DLS*, ACM, 2010, pp. 73–80.
- [12] D. Ingalls, T. Kaehler, J. Maloney, S. Wallace, A. Kay, Back to the future: the story of Squeak, a practical Smalltalk written in itself, in: *Proc. OOPSLA*, ACM, 1997, p. 318–326.
- [13] R. A. Kelsey, J. A. Rees, A tractable Scheme implementation, *Lisp Symb. Comput.* 7 (1994) 315–335.
- [14] C. F. Bolz, A. Cuni, M. Fijałkowski, A. Rigo, Tracing the meta-level: PyPy’s tracing JIT compiler, in: *ICOOOLPS*, ACM, 2009, pp. 18–25.
- [15] V. Bala, E. Duesterwald, S. Banerjia, Dynamo: a transparent dynamic optimization system, *ACM SIGPLAN Notices* 35 (2000) 1–12.
- [16] A. Gal, C. W. Probst, M. Franz, HotpathVM: an effective JIT compiler for resource-constrained devices, in: *VEE*, ACM, 2006.

- [17] C. F. Bolz, A. Cuni, M. Fijałkowski, M. Leuschel, S. Pedroni, A. Rigo, Allocation removal by partial evaluation in a tracing JIT, in: Proc. PEPM, ACM, Austin, Texas, USA, 2011, pp. 43–52.
- [18] M. Bebenita, F. Brandner, M. Fahndrich, F. Logozzo, W. Schulte, N. Tillmann, H. Venter, SPUR: a trace-based JIT compiler for CIL, in: Proc. OOPSLA, ACM, 2010, pp. 708–725.
- [19] G. Bracha, D. Ungar, Mirrors: design principles for meta-level facilities of object-oriented programming languages, in: Proc. OOPSLA, ACM, 2004, pp. 331–344.
- [20] L. Diekmann, Memory Optimizations for Data Types in Dynamic Languages, Masters thesis, Heinrich-Heine-Universität Düsseldorf, 2012.
- [21] C. F. Bolz, A. Cuni, M. Fijałkowski, M. Leuschel, S. Pedroni, A. Rigo, Runtime feedback in a meta-tracing JIT for efficient dynamic languages, in: Proc. IC00OLPS, ACM, 2011, p. 9:1–9:8.
- [22] C. Chambers, D. Ungar, E. Lee, An efficient implementation of SELF a dynamically-typed object-oriented language based on prototypes, in: OOPSLA, volume 24, ACM, 1989.
- [23] K. Barrett, B. Cassels, P. Haahr, D. A. Moon, K. Playford, P. T. Withington, A monotonic superclass linearization for Dylan, SIGPLAN Not. 31 (1996) 69–82.
- [24] A. Goldberg, Smalltalk-80: The Language and its Implementation, Addison-Wesley Series in Computer Science, Addison-Wesley, 1983.
- [25] L. P. Deutsch, A. M. Schiffman, Efficient implementation of the Smalltalk-80 system, in: POPL, ACM, Salt Lake City, Utah, 1984.
- [26] R. P. Weicker, Dhrystone benchmark (Ada version 2): rationale and measurements rules, Ada Lett. IX (1989) 60–62.
- [27] K. R. Anderson, D. Rettig, Performing Lisp analysis of the FANNKUCH benchmark, SIGPLAN Lisp Pointers VII (1994) 2–12.

Appendix A. Full experimental results

In this appendix, we present the results of running the VMs from the main paper over 11 benchmarks. The Dhrystone and Richards benchmarks are as described in the main paper; the other benchmarks are derived from the Computer Language. Because of its relative paucity of libraries, few of the benchmarks have been ported to Converge 2; we therefore exclude it from this appendix. The main aim of this appendix, therefore, is to use PyPy to understand where RPython VMs fit in the performance landscape. We believe these results are the most extensive to-date in comparing such VMs.

The Benchmarks Game often has multiple versions of the same benchmark for each language. Since none of the benchmarks has threaded versions for all languages, we have used the fastest non-threaded program for each language when possible, to avoid muddying the comparison (though the C and Java knucleotide benchmarks remain threaded, as we could find no suitable replacement). We have had to slightly modify several benchmarks for portability. We have also modified some benchmarks to run more efficiently on the VMs we are using since the Benchmarks Game sometimes targets different VMs (often older versions) that have different performance characteristics. This is fully in the spirit of the Benchmarks Game, and our changes are easily compared to the originals.

The methodology for this appendix is identical to Section 7.2. As in the main paper we measure short and long runs for each benchmark. Both experiments are fully automated and can be downloaded, for reproducibility, from http://tratt.net/laurie/research/publications/files/metatracing_vms/.

Table A.1 (split over 2 pages) shows the results from our full benchmark suite. We leave the precise interpretation of the results to readers, since one thing these 200+ data points show is that, while general trends are evident, exceptions can always be found. For example: HotSpot is always the fastest VM on long benchmarks except for RegexDNA, where interpreters beat it by some margin; and Lua is generally the fastest interpreter, except for RegexDNA. We suspect that ‘performance anomalies’ such as these will always happen. Even the best VMs have occasional weak points, sometimes by design (deliberate performance trade-offs are inevitable), sometimes because certain use cases have not yet been considered or tackled. This larger set of benchmarks gives some insight into this.

	Binary Trees			Dhrystone			Fannkuch Redux		
	14	19	50000	500000	5000000	10	10	11	
C	0.189 ± 0.008	7.986 ± 0.053	0.004 ± 0.002	0.179 ± 0.010	0.163 ± 0.006	1.992 ± 0.010			
HotSpot	0.155 ± 0.014	2.423 ± 0.043	0.107 ± 0.006	0.240 ± 0.010	0.350 ± 0.008	3.448 ± 0.029			
Lua	2.240 ± 0.071	116.810 ± 3.661	0.201 ± 0.008	19.417 ± 0.474	7.683 ± 0.321	100.536 ± 2.475			
LuaJIT	0.428 ± 0.010	23.185 ± 0.053	0.014 ± 0.006	0.879 ± 0.016	0.339 ± 0.008	4.180 ± 0.010			
CPython	4.299 ± 0.029	211.453 ± 1.117	0.368 ± 0.010	35.072 ± 0.537	9.167 ± 0.237	114.001 ± 2.189			
Jython	4.173 ± 0.063	102.303 ± 1.241	1.820 ± 0.029	28.432 ± 0.466	7.776 ± 0.419	76.069 ± 4.753			
PyPy-nonopt	0.962 ± 0.010	33.858 ± 0.135	0.127 ± 0.006	5.898 ± 0.071	1.402 ± 0.022	16.989 ± 0.220			
PyPy	0.631 ± 0.008	17.851 ± 0.088	0.069 ± 0.008	1.085 ± 0.014	1.256 ± 0.024	15.239 ± 0.223			
Ruby	1.244 ± 0.020	59.834 ± 0.990	0.312 ± 0.008	29.819 ± 0.257	13.152 ± 0.200	172.098 ± 2.168			
JRuby	1.714 ± 0.053	27.971 ± 1.223	2.050 ± 0.039	10.576 ± 0.304	6.313 ± 0.127	61.934 ± 1.513			

	Fasta			KNucleotide			Mandelbrot		
	5000000	50000000	1000000	10000000	100000000	500	500	5000	
C	0.285 ± 0.008	2.809 ± 0.039	0.413 ± 0.012	3.681 ± 0.037	0.028 ± 0.006	2.336 ± 0.008			
HotSpot	0.647 ± 0.012	5.771 ± 0.071	0.541 ± 0.043	3.303 ± 0.190	0.106 ± 0.010	2.225 ± 0.012			
Lua	6.031 ± 0.169	60.202 ± 1.080	5.469 ± 0.069	53.634 ± 1.439	0.536 ± 0.006	53.145 ± 0.084			
LuaJIT	1.373 ± 0.016	13.631 ± 0.018	1.031 ± 0.024	8.977 ± 0.094	0.031 ± 0.008	2.404 ± 0.016			
CPython	7.287 ± 0.206	71.737 ± 0.804	6.334 ± 0.063	63.405 ± 0.535	1.321 ± 0.261	134.419 ± 32.347			
Jython	18.771 ± 0.392	170.831 ± 8.967	9.443 ± 0.221	81.955 ± 1.548	4.734 ± 0.031	354.069 ± 2.885			
PyPy-nonopt	1.776 ± 0.055	16.744 ± 0.506	4.180 ± 0.041	40.879 ± 0.263	0.243 ± 0.006	19.767 ± 0.245			
PyPy	1.524 ± 0.073	14.149 ± 0.270	4.173 ± 0.037	40.863 ± 0.325	0.176 ± 0.006	13.295 ± 0.029			
Ruby	14.779 ± 0.153	147.719 ± 2.274	14.759 ± 0.378	145.703 ± 2.038	2.564 ± 0.037	255.330 ± 2.926			
JRuby	40.169 ± 1.911	388.961 ± 27.243	10.158 ± 0.212	83.895 ± 2.279	1.812 ± 0.029	45.741 ± 1.394			

	NBody			RegexDNA			RevComp		
	2500000	25000000	10000000	1000000	10000000	10000000	1000000	10000000	10000000
C	0.402 ± 0.008	3.983 ± 0.037	2.593 ± 0.031	25.899 ± 0.398	0.030 ± 0.398	0.030 ± 0.004	0.252 ± 0.010		
HotSpot	0.472 ± 0.012	3.949 ± 0.022	2.881 ± 0.410	24.922 ± 0.155	0.141 ± 0.008	0.674 ± 0.012			
Lua	14.398 ± 0.439	143.890 ± 3.653	4.959 ± 0.027	49.034 ± 0.267	0.540 ± 0.010	5.273 ± 0.104			
LuaJIT	0.652 ± 0.018	6.461 ± 0.137	4.867 ± 0.012	53.849 ± 0.110	0.188 ± 0.006	1.766 ± 0.020			
CPython	36.003 ± 0.917	359.641 ± 7.556	2.672 ± 0.043	26.484 ± 0.116	0.124 ± 0.006	1.061 ± 0.014			
Jython	42.530 ± 2.522	413.519 ± 3.847	13.022 ± 0.314	118.544 ± 2.867	1.808 ± 0.088	7.029 ± 0.290			
PyPy-nonopt	3.675 ± 0.051	36.296 ± 0.400	1.262 ± 0.014	11.348 ± 0.098	0.254 ± 0.014	2.338 ± 0.069			
PyPy	3.662 ± 0.092	36.285 ± 1.015	1.269 ± 0.022	11.482 ± 0.188	0.250 ± 0.010	2.305 ± 0.065			
Ruby	43.046 ± 0.498	431.491 ± 5.523	7.644 ± 0.076	76.401 ± 0.980	0.192 ± 0.010	1.770 ± 0.055			
JRuby	13.056 ± 1.254	114.922 ± 1.333	10.072 ± 0.688	87.696 ± 7.779	1.944 ± 0.033	4.974 ± 0.194			

	Richards			Spectral Norm		
	10	100	500	5000		
C	0.012 ± 0.006	0.079 ± 0.006	0.011 ± 0.006	1.908 ± 0.024		
HotSpot	0.109 ± 0.010	0.169 ± 0.014	0.171 ± 0.037	7.210 ± 0.010		
Lua	0.665 ± 0.024	6.574 ± 0.139	1.488 ± 0.020	147.840 ± 2.297		
LuaJIT	0.085 ± 0.006	0.763 ± 0.010	0.067 ± 0.008	6.096 ± 0.010		
CPython	1.585 ± 0.022	15.698 ± 0.227	4.766 ± 0.073	490.081 ± 2.693		
Jython	2.820 ± 0.069	13.870 ± 0.345	4.333 ± 0.100	291.221 ± 39.409		
PyPy-nonopt	0.515 ± 0.010	2.839 ± 0.016	0.100 ± 0.006	6.137 ± 0.020		
PyPy	0.267 ± 0.006	0.544 ± 0.008	0.099 ± 0.006	6.134 ± 0.010		
Ruby	0.793 ± 0.018	7.159 ± 0.061	2.897 ± 0.029	285.637 ± 2.234		
JRuby	2.130 ± 0.025	3.640 ± 0.053	3.275 ± 0.041	158.899 ± 2.914		

Table A.1: Full benchmark results.