

# Optimizing Lua using run-time type specialization

Michael Schröder

Bachelor's Thesis, March 2012

Like other dynamically typed languages, Lua spends a significant amount of execution time on type checks. Yet most programs, even if they are written in a dynamic language, are actually overwhelmingly monomorphically typed. To remove this unnecessary type-checking overhead, we implement a portable optimization scheme that rewrites virtual machine instructions at run-time based on the types of their operands. While not consistent across all platforms, we achieve average speed-ups of 1.2x on Intel, with a threaded variant of our VM showing improvements in the 1.5x to 2.4x range.

## 1 Introduction

Lua is a powerful, fast, lightweight, embeddable scripting language. [It] combines simple procedural syntax with powerful data description constructs based on associative arrays and extensible semantics. Lua is dynamically typed, runs by interpreting bytecode for a register-based virtual machine, and has automatic memory management with incremental garbage collection, making it ideal for configuration, scripting, and rapid prototyping.

<http://lua.org/about.html>

Like the language it interprets, the Lua virtual machine is small and clean, implemented in just under 15 KLoC. This makes it an excellent playground for interpreter optimization techniques. But the Lua VM is already pretty fast, compared to VMs of similar languages. Can we make it even faster?<sup>1</sup>

Usually, the first thing one attacks when trying to optimize an interpreter is its dispatch loop. The high number of branch mispredictions caused by naive dispatch implementations can have a major impact on performance [EG03b, EG03a]. However, this is only true for so called *low abstraction level interpreters*. The situation is quite different when looking at interpreters for *high-level* dynamic languages, such as Lua or Python. As Brunthaler has shown [Bru09], high abstraction level interpreters actually spend only a negligible amount of their total run-time doing instruction dispatch, because the amount of work necessary to execute each operation is comparatively high.

The obvious optimization approach, therefore, is to reduce the amount of work per operation. In any dynamically typed language, a significant amount of time is spent on type checks. When the types of values are only knowable at runtime, even a seemingly simple operation like multiplying two numbers takes up quite a lot of CPU cycles. Add operator overloading into the mix (as is possible in Lua via the mechanism of metatables) and things suddenly become rather expensive.

Yet even though the possibility of dynamic typing exists, *most "dynamically typed" programs actually aren't*. In the overwhelming majority of cases, variables will stay monomorphic throughout the

<sup>1</sup> Note that we will focus on *purely portable* optimizations, as is in the spirit of Lua. While just-in-time compilers can achieve speed-ups an order of magnitude higher compared to pure interpretation, they obviously trade speed for portability. Even if a platform is technically supported by a JIT compiler, there may be political or security-related reasons prohibiting the execution of self-modifying code (e.g. sandboxing on mobile devices). And as they are usually highly complex pieces of software, maintaining and extending JIT compilers is a decidedly non-trivial task. Besides, there already exists a JIT compiler for Lua, called LuaJIT (<http://luajit.org>), which incidentally is considered to be one of the fastest dynamic language implementations.

lifetime of the program.<sup>2</sup> This assumption provides the basis of our optimization scheme.

## 2 Run-time type specialization

The idea is this: when an instruction is dispatched for the first time, it is *specialized* according to the types of its operands, i.e. its bytecode is rewritten so that the type knowledge is now inherent in the opcode of the instruction, eliminating the need for type checks when executing the operation.<sup>3</sup> To ensure that this is safe, we have to *guard* (that is, add type checks) to any instruction that could change the type of one of the operands of the instruction we just specialized.

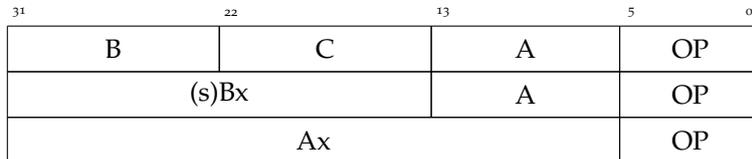
Should one of the guards fail, we will *despecialize* all instructions that are depending on that guard. For simplicity's as well as performance's sake, we adopt a very black-and-white view of the world: types are either monomorphic (in the vast majority of cases) or highly polymorphic, meaning that we will despecialize at the first sign of trouble and never respecialize again.

Put another way: we are removing type checks from *loads* of values and adding them to *stores* of values. We hope that in the end this will eliminate more type checks than it introduces.<sup>4</sup>

## 3 Prerequisites

### 3.1 Bytecode

The Lua VM is a register machine, and as such its bytecode is actually more of a *wordcode*: each instruction is exactly 32 bits long and includes an opcode and up to three operands. Operands are most commonly registers (viz. indexes into the global Lua stack, offset by the base of the function) or constants (indexes into the function's array of constant values). Depending on the operation, the bytecode is internally partitioned in one of three ways:

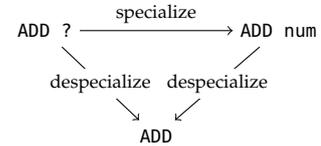


With 6 bits available for the opcode, it is clear that there can only be 64 different instructions. 40 of those are actually used by the vanilla VM. This does not leave us with nearly enough space to add all the specialized and guarded instruction variants we need, which will be just shy of 200.

One possible solution to this problem is a variable-sized bytecode, where we can have as many opcode bits as we need and fetch additional operands on demand. This would be a major change to

<sup>2</sup> An empirical study on a large number of programs written in the dynamically typed Icon language revealed a type consistency of 60% to 100%, with an average of about 80% [WG96].

<sup>3</sup> This is somewhat similar to the *quicken* technique used by the JVM and utilized by Brunthaler to implement inline caching and other optimizations for Python [Bru10a, Bru10b, Bru11].



<sup>4</sup> To see how this might work, consider that not every guarded instruction really needs to perform a type check, as the result type of an operation is often dependent on the type of its input. This will be demonstrated in greater detail in section 4.

Figure 1: Bytecode layout in the vanilla VM

the VM however, and has its own complex performance implications (cf. [OKNo2]).

Instead, we keep it simple: using the same fixed-size 32 bit instruction format, we only change the internal partitioning, to four fields of one byte each:



Figure 2: Bytecode layout in the special VM

An 8 bit opcode field gives us enough space for all the instructions we need, with room to spare. Alas, there is a tradeoff: by reducing the size of operands B and C, we impose some additional limitations on our modified VM:

	vanilla	special
possible opcodes	64	256
stack slots per function	250	120
constants per function	~67 million	~16.7 million
maximum jump offset	±131,071	±32,767

Table 1: Limits due to bytecode layout

Note that the highest-order bit of instruction fields B and C is used to differentiate between registers and constants. Thus the maximum size of the actual value in those fields is only half the size indicated by the total field length, leading to the given stack slot limits.

None of these should have any real-world impact, however. Halving the maximum number of stack slots, for example, while effectively halving the maximum number of local variables that can be declared in a single function, still leaves us with a possible 120 named variables *per function* (including function arguments and minus a few registers needed to hold temporary results). This is a limit that hopefully no sane programmer would ever come close to.

The fact that we have not run into any problems during benchmarking or when running the Lua test suite<sup>5</sup>, strongly indicates that the number of programs that can be run on the modified VM compared to vanilla Lua is not restricted in any meaningful way.

<sup>5</sup> <http://www.lua.org/tests/5.2>

### 3.2 Register scope information

The scope of a register *at a given point in time* (i.e. at some location of the program counter), is the range of instructions within which that register contains only a single semantic entity. There are two kinds of those entities:

*Local variables* All named variables, including function arguments.

The scope of a local variable ranges from the first use of that variable to the last. Note that within this range there can be multiple stores and loads to and from the register.

*Temporary variables* These arise when the result of one operation is immediately used as input to another operation. Every temporary scope begins with a register store and ends when the register is first loaded from again.

It should be clear that within a function a single register can have multiple local and temporary scopes (because it can contain different semantic entities at different points in time), but that those scopes cannot overlap.<sup>6</sup> Gathering scope information for all registers of a function can be done at compile time and is possible with only a relatively small amount of changes to the bytecode compiler: Whenever an instruction is emitted by the code generator, we decide for each of the registers used by this instruction if we should begin a new scope for the register or if we should extend the current scope to include the newly emitted instruction. This decision depends only on whether the register belongs to a local variable or not and if this particular use of the register is a load or a store. There is no need for any kind of advanced data-flow analysis.

<sup>6</sup> Although the borders of two scopes can, and quite often do, fall “within” the same instruction, e.g. `ADD 3 0 3` could end one temporary scope of register 3 and begin another one.

Observe that *each register has to be stored to at least once within a scope before it is loaded from for the first time within that same scope*. In other words: any time an instruction that loads a variable is executed, an instruction that has stored that same variable *must* have been executed recently. This is what allows us to safely transfer type checks from loads to stores.

There is a small problem, though: this guarantee does not entirely hold in the case of function arguments and return values of calls. From [IdFC05]:

For function calls, Lua uses a kind of *register window*. It evaluates the call arguments in successive registers, starting with the first unused register. When it performs the call, those registers become part of the activation record of the called function, which therefore can access its parameters as regular local variables. When this function returns, those registers are put back into the activation record of the caller.

Function arguments and return values are stored to their registers *outside* of the function they are used in, which might even be entirely outside the Lua environment. Register scopes cannot extend beyond function boundaries, which means these stores would be “invisible” to us, and there would be no way to guard them.

There is a simple solution, however: we introduce a pseudo-store instruction, which can be thought of as a “type change checkpoint”. By default, this instruction does not perform any operation, it just acts as a sentinel to inform us that at this point in the execution the contents of a certain register might have changed. We name this instruction `CHKTYPE` and it takes exactly one argument, which is the register in question. `CHKTYPE` instructions are issued at the very beginning of a function, one for each function argument, and after `CALL` and `TFCALL` instructions, one for each return value.

`CHKTYPE` instructions will also be issued after `VARARG` instructions. The `VARARG` operation stores multiple registers at once, and having one `CHKTYPE` for each store greatly simplifies specialization by allowing us to guard each register individually.

## 4 Motivational example

### 4.1 Notation

Throughout the example, we will put the specialization type of an operation (if any) to the right of the operation name and the guard type (if any) to the left. Initially, all specializable operations are marked with `?`, indicating that they have yet to be executed. Various examples of the notation are given in the table below.

GETTABLE	vanilla operation / despecialized
GETTABLE ?	not yet specialized
GETTABLE str	string-specialized
num GETTABLE	number-guarded
num GETTABLE str	string-specialized and number-guarded

Table 2: Bytecode notation examples

### 4.2 Specialization

The specialization process itself is actually rather simple, and is best explained with an example. We're going to specialize the small Lua function seen in listing 1.

The function reads `n` number of lines from the standard input and appends them to the given table `a`. Tables are associative arrays and provide the sole data structuring mechanism in Lua. Here, the table is used like a normal array. The length operator (`#`) seen on line 3 returns the number of elements in the table.

Listing 1: A small Lua function

```

1 function f(a,n)
2   for i=1,n do
3     a[#a+1] = io.read()
4   end
5 end
```

Listing 2: Vanilla bytecode

```

1 LOADK    2 -1
2 MOVE    3 1
3 LOADK    4 -1
4 FORPREP  2 6
5 LEN      6 0
6 ADD      6 6 -1
7 GETTABUP 7 0 -2
8 GETTABLE 7 7 -3
9 CALL     7 1 2
10 SETTABLE 0 6 7
11 FORLOOP  2 -7
12 RETURN   0 1
```

Listing 3: Special bytecode

```

1 CHKTYPE   0
2 CHKTYPE   1
3 LOADK    num 2 -1
4 MOVE     ? 3 1
5 LOADK    num 4 -1
6 FORPREP  ? 2 7
7 LEN      ? 6 0
8 ADD      ? 6 6 -1
9 GETTABUP ? 7 0 -2
10 GETTABLE ? 7 7 -3
11 CALL     7 1 2
12 CHKTYPE   7
13 SETTABLE ? 0 6 7
14 FORLOOP  2 -8
15 RETURN   0 1
```

Given this function, the compiler of our modified VM will produce the bytecode seen in listing 3. It differs from bytecode produced by the vanilla VM (cf. listing 2) only in the addition of `CHKTYPE` instructions (at the beginning of the function and after the `CALL` on line 12), and of course in the fact that most instructions

are now of the *not-yet-specialized* variety. Note that the two `LOADK` instructions are already specialized, since the types of their constant operands (indicated here by the use of negative indexes) are known at compile time and are guaranteed not to change.

For the purposes of our example, we are skipping over the first couple of specializations, and begin with the program counter at line 7, right before `LEN ?` is executed (see listing 4).

The single operand of `LEN ?` is register 0, containing the local variable `a`, which is a table. Examining the scope of this register, we only find one store and it is the pseudo-store of the `CHKTYPE` instruction on line 1. We add a guard to `CHKTYPE` so that it becomes `tab CHKTYPE`, which makes it safe for us to specialize `LEN ?` to `LEN tab` (see listing 5).

Next, we specialize `ADD ?`. This time there are two operands to check and guard. One of those is a constant, so no guards are necessary, and the other one is the temporary variable in register 6, whose scope ranges from the result of `LEN tab` to the input of `ADD ?`. Since the output of `LEN tab` is guaranteed to be a number, it is not necessary to add a guard. Specializing `ADD ?` to `ADD num` therefore completely eliminates the type check (see listing 6).

It should now be clear that while the details differ from operation to operation, the basic procedure is always the same:

1. Examine the instruction's operands and determine their types at this point in time and if they are suitable for specialization.
2. For each operand, guard its register so that we can safely specialize. This means looking up the scope information for that register and then examining every single instruction within that scope, with the goal of finding those instructions that *store* to the register as part of their operation. How and if we add a guard depends on the instruction in question:
  - a. If the instruction already guarantees the return type we want, then no guard is necessary. This is obviously the best case, as it completely eliminates the type check.
  - b. If the instruction does *not* already have a guaranteed return type, then guarding it *is* necessary. This effectively transfers the type check.
  - c. If instruction guarantees a return type but it is not the one we want, then we have hit upon a polymorphic type and abort the specialization process. Any other guards we might have added as part of this procedure up to now are removed again and the instruction we wanted to specialize is despecialized.
3. If adding all the necessary guards was successful, specializing the instruction is now safe. There is only one additional caveat: if the instruction we want to specialize is itself guarded, we need to reconcile this guard with the return type after specialization.

Listing 4: Before specialization of `LEN`

```

1      CHKTYPE      0
2  num  CHKTYPE      1
3      LOADK      num 2 -1
4      MOVE      num 3 1
5      LOADK      num 4 -1
6      FORPREP   num 2 7
7      LEN       ?   6 0
8      ADD       ?   6 6 -1
9      GETTABUP  ?   7 0 -2
10     GETTABLE  ?   7 7 -3
11     CALL      7 1 2
12     CHKTYPE   7
13     SETTABLE  ?   0 6 7
14     FORLOOP   2 -8
15     RETURN    0 1

```

Listing 5: After specialization of `LEN`

```

1  tab  CHKTYPE      0
2  num  CHKTYPE      1
3      LOADK      num 2 -1
4      MOVE      num 3 1
5      LOADK      num 4 -1
6      FORPREP   num 2 7
7      LEN       tab 6 0
8      ADD       ?   6 6 -1
9      GETTABUP  ?   7 0 -2
10     GETTABLE  ?   7 7 -3
11     CALL      7 1 2
12     CHKTYPE   7
13     SETTABLE  ?   0 6 7
14     FORLOOP   2 -8
15     RETURN    0 1

```

Listing 6: After specialization of `ADD`

```

1  tab  CHKTYPE      0
2  num  CHKTYPE      1
3      LOADK      num 2 -1
4      MOVE      num 3 1
5      LOADK      num 4 -1
6      FORPREP   num 2 7
7      LEN       tab 6 0
8      ADD       num 6 6 -1
9      GETTABUP  ?   7 0 -2
10     GETTABLE  ?   7 7 -3
11     CALL      7 1 2
12     CHKTYPE   7
13     SETTABLE  ?   0 6 7
14     FORLOOP   2 -8
15     RETURN    0 1

```

If the new return type is the same as the guard, we can simply remove the guard and eliminate the type check completely.

If the guard and the new return type clash, we immediately start the despecialization process on the result register, but our specialized instruction will stay specialized.

#### 4. Re-dispatch the now specialized instruction.

If we continue to apply this procedure to the rest of the function, we end up with listing 7. Take note of how the type checks that would have been necessary for `SETTABLE`, `ADD` and `LEN` could be eliminated by reducing them to the single type check at `tab` `CHKTYPE` on line 1. All in all we had to add three new type checks (on lines 1, 2 and 9), but could remove about twelve type checks by specializing the seven instructions on lines 4, 6-10 and 13.

#### 4.3 Despecialization

Our example function could be specialized by making some assumptions about its future. One of those assumption is that it will always be called with a table as the first argument. What happens if this is not the case?

We know what *should* happen: the function should behave just as it would if it was never specialized. This might either mean that it should produce a runtime exception, because whatever object we gave it instead of a table does not support the same operations, or it might mean that the function should continue to append lines read from the standard input to `a`, but maybe with a different interpretation of "append", depending on however the new object overloads the standard table operations. In either case, what the function must under no circumstances do is crash of some segmentation fault because one of the specialized instructions could no longer rely on doing things without a safety net. So we have to despecialize the function before any of this can happen. Again, the basic procedure is always the same:

1. Look up the scope information for the register whose guard has failed and examine every single instruction within that scope to find those that *load* from the register as part of their operation. What happens next depends on the particulars of each instruction:
  - a. If the instruction is not itself guarded, simply remove the instruction's specialization (e.g. rewrite `LEN tab` to `LEN`).
  - b. If the instruction *does* have itself a guard, remove the instruction's specialization but keep the guard. There is no need to mess with the specializations of other instructions as long as it is possible for guards to keep them safe.
  - c. If the instruction does not itself have a guard, but did guarantee its return type by virtue of its specialization, remove

Listing 7: Fully specialized function

```

1  tab  CHKTYPE      0
2  num  CHKTYPE      1
3      LOADK    num  2 -1
4      MOVE     num  3  1
5      LOADK    num  4 -1
6      FORPREP  num  2  7
7      LEN      tab  6  0
8      ADD      num  6  6 -1
9  tab  GETTABUP  str  7  0 -2
10     GETTABLE  str  7  7 -3
11     CALL      7  1  2
12     CHKTYPE   7
13     SETTABLE  num  0  6  7
14     FORLOOP   2 -8
15     RETURN   0  1

```

the specialization and despecialize that instruction's result register as well.<sup>7</sup>

2. Remove all remaining guards of the register within the scope.

In our example, the failing guard of `tab` `CHKTYPE` on line 1 will lead to the despecialization of `LEN` `tab`, which in turn despecializes `ADD` `num`, which in turn despecializes `SETTABLE` `num`. When the process is finished, our function will look like listing 8, and will safely work with whatever type a now has. Note that not the whole function was despecialized, just those parts relating to the polymorphic type.

#### 4.4 Upvalues

One thing that we have neglected to mention until now is how we cope with local variables that have been captured in a closure, also known as upvalues. Since Lua has full closure support, assigning to an upvalue immediately assigns to the captured local variable in the enclosing function. This means that we need to be able to despecialize within the original scope of this local variable (in the enclosing function) when changing the type of the corresponding upvalue (in the enclosed function). In essence, we need to be able to despecialize “through” upvalues, i.e. across function boundaries.

Our implementation does this and goes a step further by also allowing *specialization* to occur through upvalues, which is possible by propagating guards across the function hierarchy.<sup>8</sup>

- When specializing an instruction that has an upvalue operand (`GETUPVAL`, `GETTABUP` or `SETTABUP`), we add guards to all uses of the corresponding local variable of that upvalue in the enclosing function.
- When adding guards to a local variable, we recursively add guards to all `SETUPVAL` instructions in enclosed functions that store that local via the upvalue.
- When the type check of a guarded `SETUPVAL` fails, it despecializes the local variable in the enclosing function.
- When despecializing a local variable, we recursively despecialize all `GETUPVAL`, `GETTABUP` or `SETTABUP` instructions in enclosed functions that load that local via an upvalue.

## 5 Experimental Evaluation

### 5.1 Methodology

The sixteen benchmarks in our test suite were taken from the following sources:

- The Computer Language Benchmarks Game, a set of micro-benchmarks commonly used when comparing scripting language implementations.

<sup>7</sup> Care must be taken not to get stuck in an infinite despecialization loop by remembering which registers have already been visited.

Listing 8: After despecialization of a

```

1      CHKTYPE      0
2  num  CHKTYPE      1
3      LOADK      num 2 -1
4      MOVE      num 3 1
5      LOADK      num 4 -1
6      FORPREP   num 2 7
7      LEN       6 0
8      ADD       6 6 -1
9  tab  GETTABUP   str 7 0 -2
10     GETTABLE   str 7 7 -3
11     CALL      7 1 2
12     CHKTYPE    7
13     SETTABLE   0 6 7
14     FORLOOP    2 -8
15     RETURN     0 1

```

<sup>8</sup> Since not all of the information needed for this was available in the vanilla VM, we had to modify the virtual machine a little further so that the scopes of upvalues and their relation to parallel upvalues and enclosing functions is now fully collected at compile time and accessible to us during specialization.

- Several variations of the Richards benchmark, which simulates the task dispatcher of a simple operating system kernel. The different versions make use of different features of the Lua language, such as metatables or tail calls.
- SciMark, a popular suite of scientific and numerical benchmarks, ported to Lua by Mike Pall. We have split up the individual components and gave them fixed iteration counts so as not to get an auto-scaled score.

Both the vanilla and the special VM were compiled with `-O2 -fomit-frame-pointer`.<sup>9</sup> We measured user CPU times using the built-in `time` shell command on a variety of different platforms. For greater accuracy, benchmark inputs were chosen to produce run times around the mid double-digits whenever possible. The best results of three consecutive runs were compared.

## 5.2 Results

Based on other purely interpretative efforts to reduce type-checking overhead in virtual machines, we expected to see relative speed-ups of at least 30% [WGM10, Bru09]. The nearest we come to this is on the Intel platform, where we could achieve a 20% speed-up on average and a maximum speed-up of nearly 60%. Still, the results are somewhat underwhelming, especially when we look at the other platforms where performance was significantly worse, with not even a 10% average speed-up for two of those.

	min	avg	max
Intel Core 2 Duo SL9600	0.96	1.21	1.59
AMD Athlon 64 X2 4600+	0.96	1.08	1.22
PowerPC 970	0.95	1.11	1.26
Feroceon 88FR131 (ARM)	0.95	1.07	1.17

We are not quite sure why one of the platforms has fared so much better than the other three, but our first guess is that the increased size of the modified VM's dispatch loop (199 instructions vs. the original 40) has a negative impact on the processors' instruction cache behaviors. Differences in branch prediction accuracy might also be significant.

	L1 data	L1 instruction	L2
Intel Core 2 Duo SL9600	32 KB	32 KB	6 MB
AMD Athlon 64 X2 4600+	64 KB	64 KB	512 KB
PowerPC 970	64 KB	32 KB	512 KB
Feroceon 88FR131 (ARM)	16 KB	16 KB	256 KB

<http://lua-users.org/lists/luail/2011-04/msg00609.html>

For more information see [dQ09]

<http://luajit.org/performance.html>

<sup>9</sup> Apart from `-O2` being the default setting in Lua's makefiles, we found the results obtained using higher optimization levels to be very erratic, showing small improvements in a few cases, but exhibiting worse or unchanged performance most of the time. We attribute this to more aggressive inlining, but the full chain of cause-and-effect eludes us.

Table 3: Relative speed-ups on different platforms

For a more detailed breakdown, see Appendixes A & B

Table 4: Comparison of CPU cache sizes

### 5.3 Threading and choice of compiler

To be fully compatible with ANSI C, the Lua VM implements switch-based dispatch. We stated before that this is generally not ideal in terms of performance, but that the impact for a high level VM such as Lua should be relatively minor. This assumption turned out to be wrong: after implementing a very simple version of indirect threading,<sup>10</sup> we found all around performance gains on the order of an additional 20% on the Intel platform, before and after type specialization, compared to switch dispatch.<sup>11</sup> This goes contrary to [Bruo9] and shows that, at least for *some* high abstraction level interpreters, dispatch overhead *can* play a significant part in overall run time.

We additionally found that choice of compiler should not be underestimated. When we re-ran all benchmarks on the Intel platform for each of the modified VMs using different compilers, even though gcc proved the fastest choice for the vanilla VM, the specialized version using threading was an average 10% faster when compiled with clang, with an outlying peak speed-up of 2.45x on the mandelbrot benchmark. This goes to show that truly portable optimizations are indeed very hard to achieve, as choosing the right compiler and settings for the target platform can seriously influence the success of other optimizations.

<sup>10</sup> <http://lua-users.org/lists/lua-l/2010-11/msg00436.html>

<sup>11</sup> This is especially interesting when compared to the approach taken by Williams, McCandless and Gregg [WMG10]. With the same goal of reducing type checks, they replaced the default Lua dispatch loop wholesale with a graph-based “dynamic intermediate representation”, which models control flow and type changes along its edges. Specialized nodes allow for efficient execution of operations based on the type knowledge inherent in the structure of the graph, similar to our specialized instructions. Dispatch along this graph also reduces the number of branch mispredictions in a similar way to threading. They achieved speed-ups of 1.3x on average, with peaks at 2x, though at the cost of using significantly more memory.

While the dynamic graph allows for further, more advanced optimizations and can be used as the basis to implement a full-on JIT compiler, our simpler scheme achieves comparable performance when combined with threading.

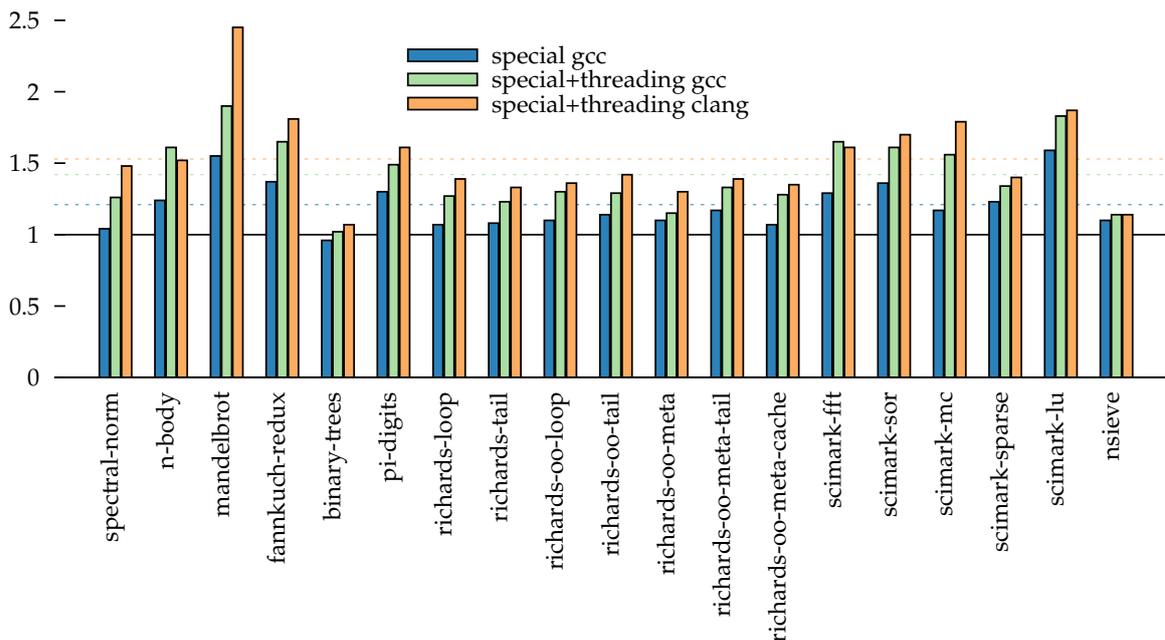


Figure 3: Relative speed-ups on the Intel platform for different optimizations and choices of compiler

## 6 Conclusion

We have presented a way to remove type-checking overhead in the Lua virtual machine by rewriting bytecode instructions at runtime. Taking advantage of the fact that most dynamically typed programs are actually highly monomorphic, we *specialize* instructions according to the types of their arguments when they are first executed, *guard* instructions on which the type assumptions made during specialization depend and *despecialize* instructions should the respective guards fail.

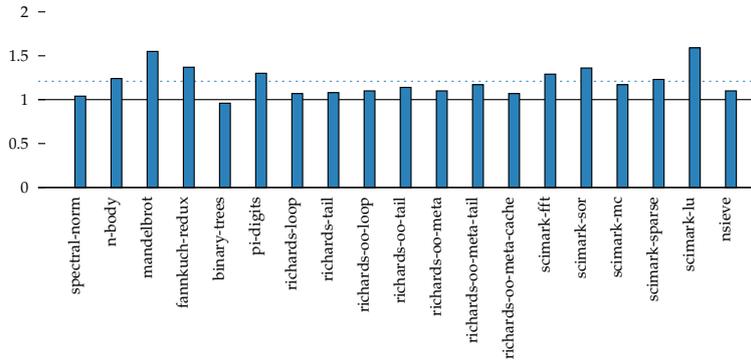
While this approach has shown promise on at least one platform, with an average speed-up of about 1.5x when combined with threaded dispatch, the results were less encouraging on other platforms. We theorized that this was most likely due to differences in CPU instruction cache and branch prediction behavior.

Adding about 1700 lines of code to the virtual machine, our implementation is still comparatively simple and straightforward. It will hopefully encourage other explorations of novel interpreter optimization techniques for the Lua language.

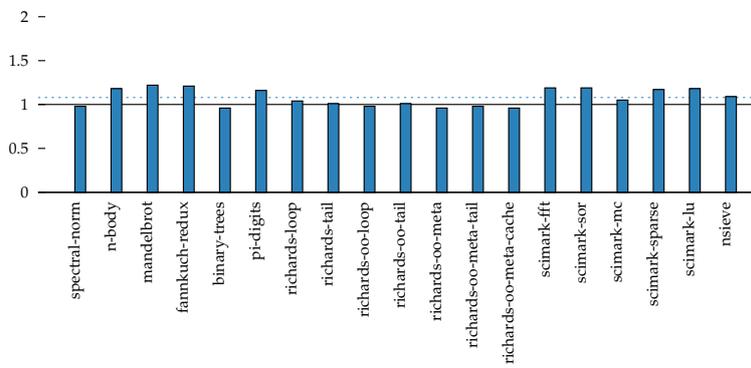
## References

- [Bru09] Stefan Brunthaler. Optimizing high abstraction-level interpreters. In *Proceedings of the 26th Annual Workshop of the GI-FG 2.1.4 "Programmiersprachen und Rechenkonzepte"* (Physikzentrum Bad Honnef, Germany, May 4-6, 2009), Institut für Informatik und Praktische Mathematik, Christian-Albrechts-Universität Kiel, Germany, Bericht Nr. 0915 (2009), pages 100–111, 2009.
- [Bru10a] Stefan Brunthaler. Efficient interpretation using quickening. In *Proceedings of the 6th Symposium on Dynamic Languages, Reno, Nevada, US, October 18, 2010 (DLS '10)*, pages 1–14, New York, NY, USA, 2010. ACM Press.
- [Bru10b] Stefan Brunthaler. Inline caching meets quickening. In *Proceedings of the 24th European Conference on Object-Oriented Programming, Maribor, Slovenia, June 21-25, 2010 (ECOOP '10)*, volume 6183 of *Lecture Notes in Computer Science*, pages 429–451. Springer, 2010.
- [Bru11] Stefan Brunthaler. *Purely Interpretative Optimizations*. PhD thesis, Vienna University of Technology, 2011.
- [dQ09] Fabio Mascarenhas de Queiroz. *Optimized Compilation of a Dynamic Language to a Managed Runtime Environment*. PhD thesis, PUC-Rio, September 2009.
- [EG03a] M Anton Ertl and David Gregg. Optimizing indirect branch prediction accuracy in virtual machine interpreters. In *SIGPLAN Conference on Programming Language Design and Implementation (PLDI'03)*, 2003.
- [EG03b] M. Anton Ertl and David Gregg. The structure and performance of efficient interpreters. *Journal of Instruction-Level Parallelism*, 5, 2003.
- [IdFC05] Roberto Ierusalimschy, Luiz Henrique de Figueiredo, and Waldemar Celes. The implementation of lua 5.0. *Journal of Universal Computer Science*, 11(7):1159–1176, July 2005.
- [OKNo2] Kazunori Ogata, Hideaki Komatsu, and Toshio Nakatani. Bytecode fetch optimization for a java interpreter. In *Proceedings of the 10th international conference on Architectural support for programming languages and operating systems, ASPLOS-X*, pages 58–67, New York, NY, USA, 2002. ACM.
- [WG96] Kenneth Walker and Ralph E. Griswold. Type inference in the icon programming language. Technical Report 93-32a, Department of Computer Science, University of Arizona, 1996.
- [WMG10] Kevin Williams, Jason McCandless, and David Gregg. Dynamic interpretation for dynamic scripting languages. In *Proceedings of the 8th annual IEEE/ACM international symposium on Code generation and optimization, CGO '10*, pages 278–287, New York, New York, USA, 2010. ACM Press.

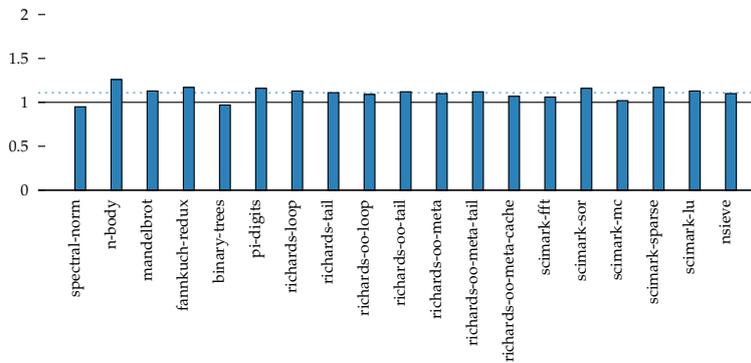
## A Relative speed-ups compared to vanilla VM



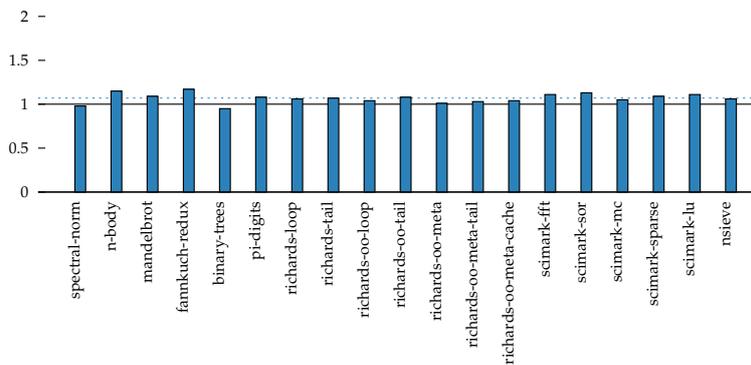
Intel Core 2 Duo SL9600 with 2.13 GHz, running OSX 10.7.3 and gcc 4.2.1



AMD Athlon 64 X2 6400+ with 2.4 GHz, running Debian 4.1.1-21 and gcc 4.1.2



PowerPC 970 2 GHz, running Debian 4.4.5-8 and gcc 4.4.5



Feroceon 88FR131 1.2 GHz, an embedded ARM device running Debian 4.4.5-2 and gcc 4.4.5

## B Detailed benchmark results

	normal		threaded		
	vanilla	special	vanilla	special	
spectral-norm	72.644	70.040	63.614	57.797	Intel Core 2 Duo SL9600 with 2.13 GHz, running OSX 10.7.3 and gcc 4.2.1
n-body	52.713	42.586	48.186	32.813	
mandelbrot	84.755	54.530	59.252	44.673	
fannkuch-redux	160.780	117.158	137.943	97.528	
binary-trees	26.628	27.802	26.205	26.051	
pi-digits	53.015	40.816	46.069	35.575	
richards-loop	52.079	48.664	46.583	40.944	
richards-tail	50.672	46.907	44.598	41.120	
richards-oo-loop	53.035	48.411	46.851	40.646	
richards-oo-tail	52.153	45.783	52.239	40.340	
richards-oo-meta	56.167	51.103	49.174	48.910	
richards-oo-meta-tail	56.274	47.985	49.582	42.464	
richards-oo-meta-cache	55.759	52.315	50.383	43.716	
scimark-fft	32.130	24.846	27.944	19.428	
scimark-sor	28.816	21.231	26.495	17.916	
scimark-mc	37.955	32.539	36.495	24.374	
scimark-sparse	39.101	31.713	36.431	29.094	
scimark-lu	46.656	29.283	34.597	25.468	
nsieve	26.297	23.812	25.562	23.038	

	normal		threaded		
	vanilla	special	vanilla	special	
spectral-norm	79.102	70.584	56.225	49.245	Intel Core 2 Duo SL9600 with 2.13 GHz, running OSX 10.7.3 and Apple clang 3.1
n-body	53.554	40.956	46.087	34.595	
mandelbrot	81.346	52.268	49.870	34.538	
fannkuch-redux	161.491	113.476	142.925	88.653	
binary-trees	26.761	26.103	25.557	24.810	
pi-digits	55.458	39.040	42.625	32.996	
richards-loop	52.829	47.525	42.739	37.429	
richards-tail	51.020	47.098	43.199	38.105	
richards-oo-loop	50.741	48.797	45.186	39.049	
richards-oo-tail	51.212	46.116	44.526	36.704	
richards-oo-meta	55.146	54.309	47.607	43.233	
richards-oo-meta-tail	55.888	50.020	48.313	40.477	
richards-oo-meta-cache	55.822	52.331	47.515	41.447	
scimark-fft	30.114	23.327	27.469	19.987	
scimark-sor	29.574	19.940	25.052	16.944	
scimark-mc	35.597	26.606	30.815	21.200	
scimark-sparse	39.459	31.499	35.104	28.003	
scimark-lu	44.674	27.959	32.581	24.909	
nsieve	26.495	23.763	24.792	23.012	

	normal		threaded	
	vanilla	special	vanilla	special
spectral-norm	114.427	116.951	106.927	103.318
n-body	63.000	53.347	54.639	42.891
mandelbrot	116.287	95.578	114.827	88.526
fannkuch-redux	194.420	160.690	176.287	140.189
binary-trees	32.894	34.430	31.510	31.882
pi-digits	76.253	65.924	61.800	49.015
richards-loop	58.012	55.795	53.851	49.463
richards-tail	57.088	56.656	54.015	50.459
richards-oo-loop	55.383	56.400	55.003	51.331
richards-oo-tail	53.707	53.151	53.507	48.663
richards-oo-meta	60.620	62.984	60.152	57.420
richards-oo-meta-tail	58.484	59.716	69.392	55.559
richards-oo-meta-cache	59.804	62.088	58.860	55.951
scimark-fft	33.138	27.826	32.646	26.430
scimark-sor	34.074	28.558	32.614	26.150
scimark-mc	45.231	43.099	40.823	35.826
scimark-sparse	46.563	39.690	42.447	35.410
scimark-lu	51.951	44.059	44.139	35.714
nsieve	35.402	32.458	30.030	28.502

AMD Athlon 64 X2 6400+ with 2.4 GHz, running Debian 4.1.1-21 and gcc 4.1.2

	normal		threaded	
	vanilla	special	vanilla	special
spectral-norm	17.86	18.80	15.07	15.20
n-body	9.26	7.35	8.19	6.62
mandelbrot	14.29	12.64	12.38	10.03
fannkuch-redux	28.52	24.44	25.37	22.11
binary-trees	13.46	13.86	12.74	12.69
pi-digits	17.96	15.54	15.98	14.11
richards-loop	13.53	11.94	11.95	11.09
richards-tail	12.96	11.69	11.68	11.41
richards-oo-loop	13.53	12.37	11.43	11.73
richards-oo-tail	13.14	11.71	11.26	11.17
richards-oo-meta	15.35	14.01	13.42	13.48
richards-oo-meta-tail	14.96	13.30	13.23	12.94
richards-oo-meta-cache	14.70	13.68	12.63	13.21
scimark-fft	11.73	11.03	10.25	8.99
scimark-sor	12.42	10.68	11.29	9.85
scimark-mc	11.16	10.93	9.78	9.33
scimark-sparse	9.12	7.81	8.11	6.90
scimark-lu	10.01	8.85	9.08	8.26
nsieve	16.73	15.25	16.32	14.82

PowerPC 970 2 GHz, running Debian 4.4.5-8 and gcc 4.4.5

	normal		threaded	
	vanilla	special	vanilla	special
spectral-norm	8.81	8.96	9.24	10.04
n-body	10.47	9.11	10.97	9.61
mandelbrot	8.65	7.93	8.83	8.36
fannkuch-redux	8.60	7.38	8.68	7.78
binary-trees	14.40	15.19	14.96	16.25
pi-digits	12.85	11.88	13.03	12.43
richards-loop	14.21	13.38	15.49	14.77
richards-tail	16.19	15.20	17.69	16.71
richards-oo-loop	15.39	14.77	16.30	16.42
richards-oo-tail	15.26	14.13	16.27	15.93
richards-oo-meta	17.60	17.47	18.47	19.11
richards-oo-meta-tail	17.35	16.77	18.28	18.24
richards-oo-meta-cache	17.40	16.73	18.38	18.76
scimark-fft	28.03	25.30	28.74	25.90
scimark-sor	27.90	24.70	28.53	25.76
scimark-mc	27.13	25.80	28.31	28.65
scimark-sparse	42.86	39.33	43.32	41.13
scimark-lu	48.07	43.18	48.39	43.39
nsieve	13.49	12.73	13.60	12.91

Feroceon 88FR131 1.2 GHz, an embedded ARM device running Debian 4.4-5-2 and gcc 4.4.5