Correctness of Speculative Optimizations with Dynamic Deoptimization

OLIVIER FLÜCKIGER, Northeastern University, USA GABRIEL SCHERER, Northeastern University, USA and INRIA, France MING-HO YEE, AVIRAL GOEL, and AMAL AHMED, Northeastern University, USA JAN VITEK, Northeastern University, USA and CVUT, Czech Republic

High-performance dynamic language implementations make heavy use of speculative optimizations to achieve speeds close to statically compiled languages. These optimizations are typically performed by a just-in-time compiler that generates code under a set of assumptions about the state of the program and its environment. In certain cases, a program may execute code compiled under assumptions that are no longer valid. The implementation must then deoptimize the program on-the-fly; this entails finding semantically equivalent code that does not rely on invalid assumptions, translating program state to that expected by the target code, and transferring control. This paper looks at the interaction between optimization and deoptimization, and shows that reasoning about speculation is surprisingly easy when assumptions are made explicit in the program representation. This insight is demonstrated on a compiler intermediate representation, named sourir, modeled after the high-level representation for a dynamic language. Traditional compiler optimizations such as constant folding, unreachable code elimination, and function inlining are shown to be correct in the presence of assumptions. Furthermore, the paper establishes the correctness of compiler transformations specific to deoptimization: namely unrestricted deoptimization, predicate hoisting, and assume composition.

CCS Concepts: • Software and its engineering → Just-in-time compilers;

Additional Key Words and Phrases: Speculative optimization, dynamic deoptimization, on-stack-replacement

ACM Reference Format:

Olivier Flückiger, Gabriel Scherer, Ming-Ho Yee, Aviral Goel, Amal Ahmed, and Jan Vitek. 2018. Correctness of Speculative Optimizations with Dynamic Deoptimization. *Proc. ACM Program. Lang.* 2, POPL, Article 49 (January 2018), 28 pages. https://doi.org/10.1145/3158137

1 INTRODUCTION

Dynamic languages pose unique challenges to compiler writers. With features such as dynamic binding, runtime code generation, and generalized reflection, languages such as Java, C#, Python, JavaScript, R, or Lisp force implementers to postpone code generation until the last possible instant. The intuition being that just-in-time (JIT) compilation can leverage information about the program state and its environment, *e.g.*, the value of program inputs or which libraries were loaded, to generate efficient code and potentially update code on-the-fly.

Many dynamic language compilers support some form of *speculative* optimization to avoid generating code for unlikely control-flow paths. In a dynamic language prevalent polymorphism causes even the simplest code to have non-trivial control flow. Consider the JavaScript snippet in Figure 1 (example from Bebenita, Brandner, Fahndrich, Logozzo, Schulte, Tillmann, and Venter [2010]). Without optimization one iteration of the loop executes 210 instructions; all arithmetic operations are dispatched and their results boxed. If the compiler is allowed to make

© 2018 Copyright held by the owner/author(s). 2475-1421/2018/1-ART49 https://doi.org/10.1145/3158137

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

for (i=0; i < a.length-1; i++) {

Fig. 1. JavaScript rotate function.

var t=a[i];a[i]=a[i+1];a[i+1]=t;

the assumption it is operating on integers, the body of the loop shrinks down to 13 instructions. As another example, most Java implementations assume that non-final methods are not overridden. Speculating on this fact allows compilers to avoid emitting dispatch code [Ishizaki, Kawahito, Yasue, Komatsu, and Nakatani 2000]. Newly loaded classes are monitored, and any time a method is overridden, the virtual machine invalidates code that contains devirtualized calls to that method. The validity of speculations is

expressed as a predicate on the program state. If some program action, like loading a new class, falsifies that predicate, the generated code must be discarded. To undo an assumption, an implementation must ensure that functions compiled under that assumption are retired. This entails replacing affected code with a version that does not depend on the invalid predicate and, if a function currently being executed is found to contain invalid code, that function needs to be replaced on-the-fly. In such a case, it is necessary to transfer control to a different version of the function, and in the process, it may be necessary to materialize portions of the state that were optimized away and perform other recovery actions. In particular, if the invalid function was inlined into another function, it is necessary to synthesize a new stack frame for the caller. This is referred to as *deoptimization*, or *on-stack-replacement*, and is found in most industrial-strength compilers.

Speculative optimization gives rise to a large and multi-dimensional design space that lies mostly unexplored. First, compiler writers must decide how to obtain information about program state. This can be done ahead-of-time by profiling, just-in-time by sampling or instrumenting code. Next, they must select what facts to record. This can range from information about the program, its class hierarchy, which packages were loaded, to information about the value of a particular mutable location in the heap. Finally, they must decide how to efficiently monitor the validity of speculations. While some points in this space have been explored empirically, existing systems have done it in an *ad hoc* manner that is often both language- and implementation-specific, and thus difficult to apply broadly.

This paper has a focused goal. We aim to demystify the interaction between compiler transformations and deoptimization. When are two versions compiled under different assumptions equivalent? How should traditional optimizations be adapted when operating on code containing deoptimization points? In what ways does deoptimization inhibit optimizations? In this work we give compiler writers the formal tools they need to reason about speculative optimizations. To do this in a way that is independent of the specific language being targeted and of implementation details relative to a particular compiler infrastructure, we have designed a high-level compiler intermediate representation (IR), named sourir, that is adequate for many dynamic languages without being tied to any one in particular.

Sourir is inspired by our work on RIR, an IR for the R language. A sourir program is made up of functions, and each function can have multiple versions. We equip the IR with a single instruction, named **assume**, specific to speculative optimization. This instruction has the role of describing what assumptions are being used to perform speculative optimization and what information must be preserved for deoptimization. It tests if those assumptions hold, and in case they do not, transfers control to another, less optimized version of the code. Reifying assumptions in the IR makes the interaction with compiler transformations explicit and simplifies reasoning. The assume instruction is more than a branch: when deoptimizing it replaces the current stack frame with a stack frame that has the variables and values expected by the target version, and, in case the function was inlined, it synthesizes missing stack frames. Furthermore, unlike a branch, its deoptimization target is not followed by the compiler during analysis and optimization. The code executed in case of deoptimization is invisible to the optimizer. This simplifies optimizations and reduces compile time as

analysis remains local to the version being optimized and the deoptimization metadata is considered to be a stand-in for the target version.

As an example consider the function from Figure 1. A possible translation to sourir is shown in Figure 2 (less relevant code elided). Vbase contains the original version. Helper functions get and store implement JavaScript (JS) array semantics, and the function add implement JS addition. Version Vnative contains only primitive sourir instructions. This version is optimized under the assumption that the variable a is an array of primitive numbers, which is represented by the first assume instruction. Further, JS arrays can be sparse and contain holes, in which case access might need to be delegated to a getter function. For this example HL denotes such a hole. The second assume instruction reifies the compiler's speculation that the array has no holes, by asserting the predicate $t \neq HL$. It also contains the associated deoptimization metadata. In case the predicate does not

```
rot()
```

```
Vnative
```

```
call type = typeof(a)
        assume type = NumArray else rot.Vbase.Lt []
  Lt
        branch i < limit Lo Lrt
  Lo
        var t = a[i]
        assume t \neq HL else rot.Vbase.Ls [i = i, j = i + 1]
        a[i] \leftarrow a[i+1]
        a[i+1] \leftarrow t
        i \leftarrow i + 1
        goto Lt
  Lrt
Vbase
  Lt
        branch i < limit Lo Lrt
        call j = add(i, 1)
  Lo
  Ls
        call t1 = get(a, i)
        call t2 = get(a, j)
        call t3 = store(a, i, t2)
        call t4 = store(a, j, t1)
        i ← j
        goto Lt
  Lrt
```

Fig. 2. Compiled function from Figure 1.

hold, we deoptimize to a related position in the base version by recreating the variables in the target scope. As can be seen in the second assume, local variables are mapped as [i = i, j = i + 1]; the current value of i is carried over into the target frame's i, whereas variable j has to be recomputed.

We prove the correctness of a selection of traditional compiler optimizations in the presence of speculation; these are constant propagation, unreachable code elimination, and function inlining. The main challenge for correctness is that the transformations operate on one version in isolation and therefore only see a subset of all possible control flows. We show how to split the work to prove correctness between the pass that establishes a version-to-version correspondence and the actual optimizations. Furthermore, we introduce and prove the correctness of three optimizations specific to speculation, namely unrestricted deoptimization, predicate hoisting, and assume composition.

Our work makes several simplifying assumptions. We use the same IR for optimized and unoptimized code. We ignore the issue of generation of versions: we study optimizations operating on a program at a certain point in time, on a set of versions created before that time. We do not model the low-level details of code generation. Correctness of runtime code generation and code modification within a JIT compiler has been addressed by Myreen [2010]. Sourir is not designed for implementation, but to give a reasoning model for existing JIT implementations. We do not intend to implement a new JIT engine. Instead, we evaluated our work by discussing it with JIT implementers; the V8 team [Chromium 2017] confirmed that intuitions and correctness arguments could be ported from sourir to their setting.

Proceedings of the ACM on Programming Languages, Vol. 2, No. POPL, Article 49. Publication date: January 2018.

2 RELATED WORK

The SELF virtual machine pioneered dynamic deoptimization [Hölzle, Chambers, and Ungar 1992]. The SELF compiler implemented many optimizations, one of which was aggressive inlining, yet the language designers wanted to give end users the illusion that they were debugging source code. They achieved this by replacing optimized code and the corresponding stack frames with non-optimized code and matching stack frames. When deoptimizing code that had been inlined, the SELF compiler synthesized stack frames. The HotSpot compiler followed from the work on SELF by introducing the idea of speculative optimizations [Paleczny, Vick, and Click 2001]. HotSpot supported very specific assumptions related to the structure of the class hierarchy and instrumented the class loader to trigger invalidation. When an invalidation occurred affected functions were rolled forward to a safe point and control was transferred from native code to the interpreter. The Jikes RVM adopted these ideas to avoid compiling uncommon code paths [Fink and Qian 2003].

One drawback of the early work was that deoptimization points were barriers around which optimizations were not allowed. Odaira and Hiraki [2005] were the first to investigate exception reordering by hoisting guards. They remarked that checking assumptions early might improve code. In Soman and Krintz [2006] the optimizer is allowed to update the deoptimization metadata. In particular they support eliding duplicate variables in the mapping and lazily reconstructing values when transferring control. This unlocks further optimizations, which were blocked in previous work. The paper also introduces the idea of being able to transfer control at any point. We support both the update of metadata and unconstrained deoptimization.

Modern virtual machines have all incorporated some degree of speculation and support for deoptimization. These include implementations of Java (HotSpot, Jikes RVM), JavaScript (WebKit Core, Chromium V8, Truffle/JS, Firefox), Ruby (Truffle/Ruby), and R (FastR), among others. Anecdotal evidence suggests that the representation adopted in this work is representative of the instructions found in the IR of production VMs: the TurboFan IR from V8 [Chromium 2017] represents assume with three distinct nodes. First a *checkpoint*, holding the deoptimization target, marks a stable point, to where execution can be rolled back. In sourir this corresponds to the original location of an assume. A *framestate* node records the layout of, and changes to, the local frame, roughly the varmap in sourir. Assumption predicates are guarded by conditional deoptimization nodes, such as *deoptimizeIf*. Graal [Duboscq, Würthinger, Stadler, Wimmer, Simon, and Mössenböck 2013] also has an explicit representation for assumptions and associated metadata as *guard* and *framestate* nodes in their high-level IR. In both cases guards are associated with the closest dominating checkpoint.

Lowering deoptimization metadata is described in Duboscq, Würthinger, and Mössenböck [2014]; Schneider and Bolz [2012]. A detailed empirical evaluation of deoptimization appears in Zheng, Bulej, and Binder [2017]. The implementation of control-flow transfer is not modeled here as it is not relevant to our results. For one particular implementation, we refer readers to D'Elia and Demetrescu [2016] which builds on LLVM. Alternatively, Wang, Lin, Blackburn, Norrish, and Hosking [2015] propose an IR that supports restricted primitives for hot-patching code in a JIT.

There is a rich literature on formalizing compiler optimizations. The CompCert project [Leroy and Blazy 2008] for example implements many optimizations, and contains detailed proof arguments for a data-flow optimization used for constant folding that is similar to ours. In fact, sourir is close to CompCert's RTL language without versions or assumptions. There are formalizations for tracing compilers [Dissegna, Logozzo, and Ranzato 2014; Guo and Palsberg 2011], but we are unaware of any other formalization effort for speculative optimizations in general. Béra, Miranda, Denker, and Ducasse [2016] present a verifier for a bytecode-to-bytecode optimizer. By symbolically executing optimized and unoptimized code, they verify that the deoptimization metadata produced by their optimizer correctly maps the symbolic values of the former to the latter at all deoptimization points.

3 SOURIR: SPECULATIVE COMPILATION UNDER ASSUMPTIONS

This section introduces our IR and its design principles. We first present the structure of programs and the assume instruction. Then, Section 3.2 and following explain how sourir maintains multiple equivalent versions of the same function, each with a different set of assumptions. This enables the speculative optimizations presented in Section 4. All concepts introduced in this section are formalized in Section 5.

3.1 Sourir in a Nutshell

Sourir is an untyped language with lexically scoped mutable variables and first-class functions. As an example the function in Figure 3 queries a number n from the user and initializes an array with values from 0 to n-1. By design, sourir is a cross between a compiler representation and a high-level language. We have equipped it with sufficient expressive power so that it is possible to write interesting programs in a style reminiscent of dynamic languages.¹ The only features that are critical to our result are *versions* and *assumptions*. Versions are the counterpart of dynamically generated code fragments. Assumptions, represented by the assume instruction, support dynamic deoptimization of speculatively compiled code. The syntax of sourir instructions is shown in Figure 4.

var n = nil read n array t[n] var k = 0 goto L1 L1 branch k < n L2 L3 L2 t[k] \leftarrow k k \leftarrow k + 1 goto L1 L3 drop k stop

Fig. 3. Example sourir code.

Sourir supports defining a local variable, removing a vari-

able from scope, variable assignment, creating arrays, array assignment, (unstructured) control flow, input and output, function calls and returns, assumptions, and terminating execution. Control-flow instructions take explicit labels, which are compiler-generated symbols but we sometimes give them meaningful names for clarity of exposition. Literals are integers, booleans, and nil. Together with variables and function references, they form simple expressions. Finally, an expression is either a simple expression or an operation: array access, array length, or primitive operation (arithmetic, comparison, and logic operation). Expressions are not nested—this is common in intermediate representations such as A-normal form [Sabry and Felleisen 1992]. We do allow bounded nesting in instructions for brevity.

A program P is a set of function declarations. The body of a function is a list of versions indexed by a version label, where each version is an instruction sequence. The first instruction sequence in the list (the *active version*) is executed when the function is called. F ranges over function names, Vover version labels, and L over instruction labels. An absolute reference to an instruction is thus a triple F.V.L. Every instruction is labeled, but for brevity we omit unused labels.

Versions model the speculative optimizations performed by the compiler. The only instruction that explicitly references versions is assume. It has the form **assume** e^* else $\xi \ \tilde{\xi}^*$ with a list of predicates (e^*) and deoptimization metadata ξ and $\tilde{\xi}^*$. When executed, assume evaluates its predicates; if they hold execution skips to the next instruction. Otherwise, deoptimization occurs according to the metadata. The format of ξ is $F.V.L \ [x_1 = e_1, \dots, x_n = e_n]$, which contains a target F.V.L and a varmap $[x_1 = e_1, \dots, x_n = e_n]$. To deoptimize, a fresh environment for the target is created according to the varmap. Each expression e_i is evaluated in the old environment and bound to x_i in the new environment. The environment specified by ξ replaces the current one. Deoptimization might also need to create additional continuations, if assume occurs in an inlined function. In this case multiple $\tilde{\xi}$ of the form $F.V.L \ x \ [x_1 = e_1, \dots, x_n = e_n]$ can be appended. Each one synthesizes a continuation

¹An implementation of sourir and the optimizations presented here is available at https://github.com/reactorlabs/sourir.

i ::= 	var x = e drop x $x \leftarrow e$ array x[e] array x = [e*]	instructions variable declaration drop a variable from scope assignment array allocation array creation	е	::= 	se x[se] length(se) primop (se*)	expression simple expression array access array length primitive operation
	$x[e_1] \leftarrow e_2$ branch $e L_1 L_2$ goto L print e read x	array assignment conditional branch unconditional branch print read	se	::= 	lit F x	simple expressions literals function reference variables
	call $x = e(e^{-1})$ return e^{-1} assume e^{*} else $\xi \ \tilde{\xi}^{*}$ stop	return assume instruction terminate execution	lit	::= 	, -1, 0, 1, nil true false	literals numbers others
	Ę Ę			ta: ex	rget and varmap tra continuation	

ξ	::=	F.V.L x VA	extra continuation
VA	::=	$[x_1 = e_1, \ldots, x_n = e_n]$	varmap

Fig. 4. The syntax of sourir.

with an environment constructed according to the varmap, a return target F.V.L, and the name xto hold the returned result-this situation and inlining are discussed in Section 4.3. The purpose of deoptimization metadata is twofold. First, it provides the necessary information for jumping to the target version. Second, its presence in the instruction stream allows the optimizer to keep the mapping between different versions up-to-date.

Example. Consider the function size in Figure 5 which computes the size of a vector x. In version Vb, x is either nil or an array with its length stored at index 0. The optimized version Vo expects that the input is never nil. Classical compiler optimizations can leverage this fact: unreachable code removal prunes the unused branch. Constant propagation replaces the use of el with its value and updates the varmap so that it restores the deleted variable upon deoptimization to the base version Vb.

```
size(x)
   Vo
           assume x \neq nil else size.Vb.L2 [el = 32, x = x]
          var I = x[0]
          return | * 32
   Vb
     L1 var el = 32
     L2 branch x = nil L3 L4
     L3 var | = x[0]
          return | * el
     L4 return 0
```

Fig. 5. Speculation on x.

3.2 Deoptimization Invariants

A version is the unit of optimization and deoptimization. Thus we expect that each function will have one original version and possibly many optimized versions. Versions are constructed such that they preserve two crucial invariants: (1) *version equivalence* and (2) *assumption transparency*. By the first invariant all versions of a function are observationally equivalent. The second invariant ensures that even if the assumption predicates *do* hold, deoptimizing to the target should be correct. Thus one could execute an op-

show(x) Vo assume x = 42 else show.Vb.L1 [x = x] print 42 Vw assume true else show.Vb.L1 [x = 42] print x Vb L1 print x

Fig. 6. The version w violates the deoptimization invariant.

timized version and its base in lockstep; at every assume the varmap provides a complete mapping from the new version to the base. This simulation relation between versions is our correctness argument. The transparency invariant allows us to add assumption predicates without fear of altering program semantics. Consider a function show in Figure 6 which prints its argument x. Version Vo respects both invariants: any value for x will result in the same behavior as the base version and deoptimizing is always possible. On the other hand, Vw, which is equivalent because it will never deoptimize, violates the second invariant: if it were to deoptimize, the value of x would be set to 42, which is almost always incorrect. We present a formal treatment of the invariants and the correctness proofs in Section 5.4 and following.

3.3 Creating Fresh Versions

We expect that versions are chained. A compiler will create a new version, say V1, from an existing version V0 by copying all instructions from the original version and chaining their deoptimization targets. The latter is done

by updating the target and varmap of assume instructions such that all targets refer to V₀ at the same label as the current instruction. As the new version starts out as a copy, the varmap is the identity function. For instance, if the target contains the variables x and y, then the varmap is [x = x, z = z]. Additional assume instructions can be added; assume instructions that bear no predicates (i.e., the predicate list is either empty or just tautologies) can be removed while preserving equivalence. As an example in Figure 7, the new version V2 is a copy of V1; the instruction at L0 was added, the instruction at L1 was updated, and the one at L2 was removed.

fun() V2 Lo assume true else fun.V1.L0 [] var x = 1**assume** *e* **else** fun.V1.L1 [x = x]L1 L2 print x + 2V1 Lo var x = 1**assume** *e* **else** fun.V0.L1 [g = x]L1 **assume true else** fun.V0.L2 [g = x, h = x + 1]L2 print x + 2V0 Lo var g = 1var h = g + 1L1 L2 print h + 1

Fig. 7. Chained assume instructions: Version 1 was created from 0, then optimized. Version 2 is a fresh copy of 1.

Updating assume instructions is not required for correctness. But the idea with a new version is that it captures a set of assumptions that can be undone independently from the previously existing assumptions. Thus, we want to be able to undo one version at a time. In an implementation, versions might, for example, correspond to optimization tiers.² This approach can lead to a cascade of deoptimizations if an inherited assumption fails; we discuss this in Section 4.6. In the following sections we

```
size(x)
Vdup
L1 assume true else size.Vb.L1 [x = x]
var el = 32
L2 assume true else size.Vb.L2 [el = el, x = x]
branch x = nil L5 L6
L5 var l = x[0]
return l * el
L6 return 0
Vb ...
```

Fig. 8. A fresh copy of the base version of size.

use the base version Vb of Figure 5 as our running example. As a first step, we generate the new version Vdup with two fresh assume instructions shown in Figure 8. Initially the predicates are true and the assume instructions never fire. Version Vb stays unchanged.

3.4 Injecting Assumptions

We advocate an approach where the compiler first injects assumption predicates, and then uses them in optimizations. In contrast, earlier work would apply an unsound optimization and then recover by adding a guard (see, for example, Duboscq et al. [2013]). While the end result is the same, the different perspective helps with reasoning about correctness. Assumptions are boolean predicates, similar to user-provided assertions. For example, to speculate on a branch target, the assumption is the branch condition or its negation. It is therefore correct for the compiler to expect that the predicate holds immediately following an assume. Injecting predicates is done after establishing the correspondence between two versions with assume instructions, as presented above. Inserting a fresh assume in a function is difficult in general, as one must determine where to transfer control to or how to reconstruct the target environment. On the other hand, it is always correct to add a predicate to an existing assume. Thanks to the assumption transparency invariant it is safe to deoptimize more often to the target. For instance, in **assume** $x \neq nil$, x > 10 else ... the predicate $x \neq nil$ was narrowed down to x > 10.

4 OPTIMIZATION WITH ASSUMPTIONS

In the previous section we introduced our approach for establishing a fresh version of a function that lends itself to speculative optimizations. Next, we introduce classical compiler optimizations that are exemplary of our approach. Then we give additional transformations for the assume in Section 4.4 and following, and conclude with a case study in Section 4.7. All transformations introduced in this section are proved correct in Section 6.

4.1 Constant Propagation

Consider a simple constant propagation pass that finds constant variables and then updates all uses. This pass maintains a map from variable names to constant expressions or *unknown*. The map is computed for every position in the instruction stream using a data-flow analysis. Following the approach by Kildall [1973], the analysis has an update function to add and remove constants to the map. For example analyzing **var** x = 2, or $x \leftarrow 2$ adds the mapping $x \rightarrow 2$. The instruction

 $^{^{2}}$ A common strategy for VMs is to have different kind of optimizing compilers with different compilation speed versus code quality trade-offs. The more a code fragment is executed, the more powerful optimizations will be applied to it.

var y = x + 1 adds $y \rightarrow 3$ to the previous map. Finally, **drop** x removes a mapping. Control-flow merges rely on a join function for intersecting two maps; mappings which agree are preserved, while others are set to *unknown*. In a second step, expressions that can be evaluated to values are replaced and unused variables are removed. No additional care needs to be taken to make this pass correct in the presence of assumptions. This is because in sourir, the expressions needed to reconstruct environments appear in the varmap of the assume and are thus visible to the constant propagation pass. Additionally, the pass can update them, for example, in **assume true else** F.V.L [x = y + z], the variables y and z are treated the same as in **call** h = foo(y + z). They can be replaced and will not artificially keep constant variables alive.

Constant propagation can become speculative. After the instruction **assume** x = 0 else ..., the variable x is 0. Therefore, $x \leftarrow 0$ is added to the state map. This is the only extension required for speculative constant propagation. As an example, in the case where we speculate on a nil check

L2 **assume** $x \neq nil$ else size.Vb.L2 [el = el, x = x] **branch** x = nil L5 L6

the map is $x \to \neg nil$ after L2. Evaluating the branch condition under this context yields $\neg nil == nil$, and a further optimization opportunity presents itself.

4.2 Unreachable Code Elimination

As shown above, an assumption coupled with constant folding leads to branches becoming deterministic. Unreachable code elimination benefits from that. We consider a two step algorithm: the first pass replaces **branch** $e \, L_1 \, L_2$ with **goto** L_1 if e is a tautology and with **goto** L_2 if it is a contradiction. The second pass removes unreachable instructions. In our running example from Figure 8, we add the predicate $x \neq nil$ to the empty assume at L2. Constant propagation shows that the branch always goes to

```
size(x)

Vpruned

L1 assume true else size.Vb.L1 [x = x]

var el = 32

L2 assume x \neq nil else size.Vb.L2 [el = el, x = x]

var l = x[0]

return l * el

Vb ...
```

Fig. 9. A speculation that the argument is not nil eliminated one of the former branches.

L5, and unreachable code elimination removes the dead statement at L6 and branch. This creates the version shown in Figure 9. Additionally, constant propagation can replace el by 32. By also replacing its mention in the varmap of the assume at L2, el becomes unused and can be removed from the optimized version. This yields version Vo in Figure 5 at the top.

4.3 Function Inlining

Function inlining is our most involved optimization, since assume instructions inherited from the inlinee need to remain correct. The inlining itself is standard. Name mangling is used to separate the caller and callee environments. As an example Figure 10 shows the inlining of size into a function main. Naïvely inlining without updating the metadata of the assume at L1 will result in an incorrect deoptimization, as execution would transfer to size.Vb.L2 with no way to return to the main function. Also, main's part of the environment is discarded in the transfer and permanently lost. The solution is to synthesize a new stack frame. As shown in the figure, the assume at in the optimized main is thus extended with main.Vb.Lret s [pl = pl, vec = vec].

Flückiger, Scherer, Yee, Goel, Ahmed, and Vitek

This creates an additional stack frame that returns to the base version of main. and stores the result in s with the entire caller portion of the environment reconstructed. It is always possible to compute the continuation, since the original call site must have a label and the scope at this label is known. Overall, after deoptimization, it appears as if version Vb of main had called version Vb of size. Note, it would erroneous to create a continuation that returns to the optimized version of the caller Vinl. If deoptimization from the inlined code occurs, it is precisely because some of its assumptions are invalid. Multiple continuations can be appended for further levels of inlining. The inlining needs to be applied bottom up: for the next level of inlining, e.g., to inline Vinl into an outer caller, renamings must also be applied to the expressions in the extra continuations, since they refer to local variables in Vinl.

main() Vinl **array** pl = [1, 2, 3, 4] **array** vec = [length(pl), pl] var s = nil**var** x = vecsize.Vb.L2 [el = 32, x = x]assume $x \neq$ nil else main.Vb.Lret s [pl = pl, vec = vec]var | = x[0] $s \leftarrow | * 32$ drop drop x goto Lret Lret print s stop Vb **array** pl = [1, 2, 3, 4] **array** vec = [length(pl), pl] **call** s = size(vec) Lret print s stop size(x)Vo L1 **assume** $x \neq$ **nil else** size.Vb.L2 [el = 32, x = x] var | = x[0]return 1 * 32 Vb ...

Fig. 10. An inlining of size into a main.

4.4 Unrestricted Deoptimization

The assume instructions are expensive: they create dependencies on live variables and are barriers for moving instructions. Hoisting a side-effecting instruction over an assume is invalid, because if we deoptimize the effect happens twice. Removing a local variable is also not possible if its value is needed to reconstruct the target environment. Thus it makes sense to insert as few assume instructions as possible. On the other hand it is desirable to be able to "deoptimize everywhere"—checking assumptions in the basic block in which they are used can avoid unnecessary deoptimization—so there is a tension between speculation and optimization. Reaching an assume marks a stable state in the execution of the program that we can fall back to, similar to a transaction. Implementations, like [Duboscq et al. 2013], separate deoptimize more freely. As long as the effects of instructions performed since the last deoptimization point are not observable, it is valid to throw away intermediate results and resume control from there. Effectively, in sourir this corresponds to moving an assume instruction forward in the instruction stream, while keeping its deoptimization target fixed. An assume can be moved over another instruction if that instruction:

- (1) has no side-effects and is not a call instruction,
- (2) does not interfere with the varmap or predicates, and
- (3) has the assume as its only predecessor instruction.

The first condition prevents side-effects from happening twice. The second condition can be enabled by copying the affected variables at the original assume instruction location (*i.e.*, taking a snapshot of the required part of the environment).³ The last condition prevents capturing traces incoming from other basic blocks where (1) and (2) do not hold for all intermediate instructions since the original location. This is not the weakest condition, but a reasonable, sufficient one. Let us consider a modified version of our running example in Figure 11 on the left. Again, we have an assume before the branch, but would like to place a guard inside one of the branches.

size(x) Vany L1 L5 L6	assume true else size.Vb.L1 $[x = x]$ var el = 32 branch x = nil L6 L5 $x \leftarrow x[0]$ return x * el 	size(x) Vany Vary var x0 = x var el = 32 branch x = nil L6 L5 L5 $x \leftarrow x[0]$ L1 assume x = 1 else size.Vb.L1 [x = return 1 * el L6	= x0]
Vb		Vb	

Fig. 11. Moving an assume forward in the instruction stream.

There is an interfering instruction at L5 that modifies x. By creating a temporary variable to hold the value of x at the original assume location, it is possible to resolve the interference. Now the assume can move inside the branch and a predicate can be added on the updated x (see right side of the figure). Note that the target is unchanged. This approach allows for the (logical) separation between the deoptimization point and the position of assumption predicates. In the transformed example a stable deoptimization point is established at the beginning of the function by storing the value of x, but then the assumption is checked only in one branch. The intermediate states are ephemeral and can be safely discarded when deoptimizing. For example the variable el is not mentioned in the varmap here, it is not captured by the assume. Instead it is recomputed by the original code at the deoptimization target size.Vb.L1. To be able to deoptimize from any position it is sufficient to have an assume after every side-effecting instruction, call, and control-flow merge.

4.5 Predicate Hoisting

Moving an assume backwards in the code would require replaying the moved-over instructions in the case of deoptimization. Hoisting **assume true else** size.Vb.L2 [el = el, . . .] above **var** el = 32 is allowed if the varmap is changed to [el = 32, . . .] to compensate for the lost definition. However this approach is tricky and does not work for instructions with multiple predecessors as it could lead to conflicting compensation code. But a simple alternative to hoisting assume is to hoist a *predicate* from one assume to a previous one. To understand why, let us decompose the approach into two steps. Given an assume at L1 that dominates a second one at L2, we copy a predicate from the latter to the former. This is valid since the assumption transparency invariant allows strengthening predicates. A data-flow analysis can determine if the copied predicate from L1 is available at L2, in which case it can be removed from the original instruction. In our running example, version Vpruned in Figure 9 has two assume instructions and one predicate. It is trivial to hoist x \neq nil, since there are no interfering instructions. This allows us to remove the assume with the larger

³In an SSA based IR this step is not necessary for SSA variables, since the captured ones are guaranteed to stay unchanged.

scope. More interestingly, in the case of a loop-invariant assumption, predicates can be hoisted out of the loop.

4.6 Assume Composition

As we have argued in Section 3.3, it is beneficial to undo as few assumptions as possible. On the other hand, deoptimizing an assumption added in an early version cascades through all the later versions. To be able to remove chained assume instructions, we show that assumptions are *composable*. If an assume in version V3 transfers control to a target V2.La, that is itself an assumption with V1.Lb as target, then we can combine the metadata to take both steps at once. By the assumption transparency invariant, the pre- and post-deoptimization states are equivalent: even if the assumptions are not the same, it is correct to conservatively trigger the second deoptimization. For example, an instruction **assume** *e* else F.V2.La [x = 1] that jumps to **assume** *e*' else F.V0.Lb [y = x] can be combined as **assume** *e*, *e*' else F.V0.Lb [y = 1]. This new unified assume skips the intermediate version V2 and goes to V0 directly. This could be an interesting approach for multi-tier JITs: after the system stabilizes, intermediate versions are rarely used and may be discarded.

4.7 Case Study

We conclude with an example. In dynamic languages code is often dispatched on runtime types. If types were known, code could be specialized, resulting in faster code with fewer checks and branches. Consider Figure 12(a) which implements a generic binary division function that expects two values and their type tags. No static information is available; the arguments could be any type. Therefore, multiple checks are needed before the division; for example the slow branch will require even more checks on the exact value of the type tag. Suppose there is profiling information that indicates numbers can be expected. The function is specialized by speculatively pruning the branches as shown in Figure 12(b). In certain cases, sourir's transformations can make it appear as though checks have been reordered. Consider a variation of the previous example, that speculates on x, but not y as shown in Figure 12(c). In this version, both checks on x are performed first and then the ones on y, whereas in the unoptimized version they are interleaved. By ruling out an exception early, it is possible to perform the checks in a more efficient order. The fully speculated on version contains only the integer division and the required assumptions (Figure 12(d)). This version has no more branches and is a candidate for inlining.

```
div(tagx, x, tagy, y)
   Vbase
                                                           assume tagx = NUM, tagy = NUM else div.Vb.L1 [...]
             branch tagx ≠ NUM Lslow L2
     11
                                                          branch x = 0 Lerror L4
             branch tagy ≠ NUM Lslow L3
     L2
                                                      L4 return y/x
     L3
             branch x = 0 Lerror L4
                                                           . . .
     L4
             return y/x
     Lslow
            . . .
                      (a)
                                                                                    (b)
     assume tagx = NUM, x \neq 0 else div.Vb.L1 [...]
     branch tagy ≠ NUM Lslow L4
                                                           assume tagx = NUM, tagy = NUM, x \neq 0 else div.Vb.L1 [...]
L4 return y/x
                                                           return y/x
                      (c)
                                                                                    (d)
```

Fig. 12. Case study.

5 SPECULATIVE COMPILATION FORMALIZED

A sourir program contains several functions, each of which can have multiple versions. This highlevel structure is described in Figure 13. The first version is considered the currently active version and will be executed by a call instruction. Each version consists of a stream of labeled instructions. We use an indentation-based syntax that directly reflects this structure and omit unreferenced instruction labels.

$P ::= \begin{bmatrix} F(x^*) \\ V \\ L & i \end{bmatrix}$	indentation-based syntax
--	--------------------------

P::= $F(x^*): D_F, \ldots$ a program is a list of named functions D_F ::= $V: I, \ldots$ a function definition is a list of versioned instruction streamsI::= $L: i, \ldots$ an instruction stream with labeled instructions

Fig. 13. Program syntax.

Besides grammatical and scoping validity, we impose the following well-formedness requirements to ease analysis and reasoning. The last instruction of each version of the main function is **stop**. Two variable declarations for the same name cannot occur in the same instruction stream. This simplifies reasoning by letting us use variable names to unambiguously track information depending on the declaration site. Different versions have separate scopes and can have names in common. If a function reference *F* is used, that function *F* must exist. Source and target of control-flow transitions must have the same set of declared variables. This eases determining the environment at any point. To jump to a label *L*, all variables not in scope at *L* must be dropped (**drop** x).

5.1 Operational Semantics: Expressions

Figure 14 gives the semantics of expressions. Evaluation *e* returns a value *v*, which may be a literal *lit*, a function, or an address *a*. Arrays are represented by addresses into heap *M*. The heap is a map from addresses to blocks of values $[v_1, ..., v_n]$. An environment *E* is a mapping from variables to values. Evaluation is defined by a relation $M E e \rightarrow v$: under *M* and environment *E*, *e* evaluates to *v*. This definition in turn relies on a relation $E se \rightarrow v$ defining evaluation of simple expressions *se*, which does not access arrays. The notation [[primop]] to denote, for each primitive operation *primop*, a partial function on values. Arithmetic operators and arithmetic comparison operators are only defined when their arguments are numbers. Equality and inequality are defined for all values. The relation $M E e \rightarrow v$, when seen as a function from *M*, *E*, *e* to *v*, is partial: it is not defined on all inputs. For example, there is no *v* such that the relation $M E x[se] \rightarrow v$ holds if E(x) is not an address *a*, if *a* is not bound in *M*, if *se* does not reduce to a number *n*, or if *n* is out of bounds.

5.2 Operational Semantics: Instructions and Programs

We define a small-step, labeled operational semantics with a notion of machine state, or configuration, that represents the dynamic state of a program being executed, and a transition relation between configurations. A configuration is a six-component tuple $\langle PILK^* ME \rangle$ described in Figure 15. Continuations *K* are tuples of the form $\langle ILxE \rangle$, storing the information needed to correctly return to a caller function. On a call **call** $x = e(e_1, ..., e_n)$, the continuation pushed on the stack contains the current instruction stream *I* (to be restored on return), the label *L* of the

v ::= values		
lit	addr ::= a	addresses
<i>F</i>	$M ::= (a \rightarrow [v_1, \ldots, v_n])^*$	heap
a	$E ::= (\mathbf{x} \to \mathbf{v})^*$	environment
[LITERAL]	[Funref]	[Lookup]
$\overline{E \ lit} \rightarrow lit$	$\overline{E F \rightharpoonup F}$	$\overline{E x \rightharpoonup E(x)}$
[SimpleExp]		[Primop]
$E se \rightarrow v$	$E se_1 \rightharpoonup v_1 \dots E se_n \rightharpoonup v_n$	
$\overline{M E se \to v}$	$\overline{M \ E \ primop(se_1, \ldots, se_n)} \to \llbracket p$	$primop]](v_1, \ldots, v_n)$
		[VecAccess]
[VecLen]	$a \stackrel{\text{def}}{=} E(x)$	$M(a) = [v_0, \ldots, v_m]$
$E se \rightarrow a M(a) = [v_1, \ldots, v_n]$	$E se \rightarrow n$	$0 \le n \le m$
$M E \text{ length}(se) \rightarrow n$	M E	$x[se] \rightarrow v_n$



next instruction after the call (the return label), the variable x to name the returned result, and environment E. For the details, see the reduction rules for **call** and **return** in Figure 17.

$$C ::= \langle PILK^* ME \rangle \begin{pmatrix} P & & \text{program} \\ I & & \text{instructions} \\ L & & \text{next label} \\ K^* & ::= & (K_1, \dots, K_n) & \text{call stack} \\ M & & \text{heap} \\ E & & \text{environment} \\ K ::= \langle ILxE \rangle \\ \text{continuation} & \begin{pmatrix} I & & \text{code of calling function} \\ L & & \text{return label} \\ x & & \text{return variable} \\ E & & \text{environment at call site} \\ \end{pmatrix}$$

Fig. 15. Abstract machine state.

The relation $C \xrightarrow{A_{\tau}} C'$ specifies that executing the next instruction may result in the configuration C'. The action A_{τ} indicates whether this reduction is observable: it is either the silent action, written τ , an I/O action read *lit* or print *lit*, or stop. We write $C \xrightarrow{T} C'$ when there are zero or more steps from C to C'. The trace T is a list of non-silent actions in the order in which they appeared. Actions are defined in Figure 16, and the full reduction relation is given in Figure 17.

Most rules get the current instruction, I(L), perform an operation, and advance to the next label, referred to by the shorthand (L+1). The read *lit* and print *lit* actions represent observable I/O operations. They are emitted by READ and PRINT in Figure 17. The action read *lit* on the **read** x transition may be any literal value. This is the only reduction rule that is non-deterministic. Note that the relation $C \longrightarrow^* C'$, containing only sequences of silent reductions, is deterministic. The



Fig. 16. Actions and traces.

stop reduction emits the stop transition, and also produces a configuration with no instructions, \emptyset . This is a technical device to ensure that the resulting configuration is stuck. A program with a silent loop has a different trace from a program that halts. Given a program *P*, let start(*P*) be its starting configuration, and reachable(*P*) be the set of configurations reachable from it; they are all the states that may be encountered during a valid run of *P*.

$$[STARTCONF]$$

$$I \stackrel{\text{def}}{=} P(\text{main, active}) \quad L \stackrel{\text{def}}{=} \text{start}(I)$$

$$start(P) \stackrel{\text{def}}{=} \langle PIL \emptyset \ \emptyset \emptyset \rangle$$
reachable(P) \stackrel{\text{def}}{=} \{C \mid \exists T, \text{ start}(P) \stackrel{T}{\longrightarrow} C\}

5.3 Equivalence of Configurations: Bisimulation

We use weak bisimulation to prove equivalence between configurations. The idea is to define, for each program transformation, a correspondence relation *R* between configurations over the source and transformed programs. We show that related configurations have the same observable behavior, and reducing them results in configurations that are themselves related. Two programs are equivalent if their starting configurations are related.

Definition 5.1 (Weak Bisimulation). Given programs P_1 and P_2 and relation R between the configurations of P_1 and P_2 , R is a weak simulation if for any related states $(C_1, C_2) \in R$ and any reduction $C_1 \xrightarrow{A_\tau} C'_1$ over P_1 , there exists a reduction $C_2 \xrightarrow{A_{\tau}} C'_2$ over P_2 such that (C'_1, C'_2) are themselves related by R. Reduction over P_2 is allowed to take zero or more steps, but not to change the trace. In other words, the diagram on the left below can always be completed into the diagram on the right.



R is a weak bisimulation if it is a weak simulation and the symmetric relation R^{-1} also is—a reduction from C_2 can be matched by C_1 . Finally, two configurations are *weakly bisimilar* if there exists a weak bisimulation *R* that relates them.

In the remainder, the adjective weak is always implied. The following result is standard, and essential to compose the correctness proof of subsequent transformation passes.

[Decl]	[Drop]
$I(L) = $ var $x = e $ $M E e \rightarrow v$	$I(L) = \operatorname{drop} x$
$\overline{\langle PILK^* ME \rangle} \xrightarrow{\tau} \langle PI(L+1)K^* ME[x \leftarrow v] \rangle$	$\langle PILK^* ME \rangle \xrightarrow{\tau} \langle PI(L+1)K^* ME \setminus \{x\} \rangle$
[Arr	ayDef] [ArrayDeci
$I(L) = \operatorname{array} x = [e_1, \ldots, e_n] M E e_1 \rightarrow v_1 \ldots M E$	$e_n \to v_n$ $I(L) = \operatorname{array} x[e] M E e \to n$
$a \text{ fresh } M' \stackrel{\text{def}}{=} M[a \leftarrow [v_1, \ldots, v_n]]$	$a \text{ fresh } M' \stackrel{\text{def}}{=} M[a \leftarrow [\mathbf{nil}_1,, \mathbf{nil}_n]]$
$\langle P I L K^* M E \rangle \xrightarrow{\tau} \langle P I (L+1) K^* M' E[x \leftarrow a] \rangle$	$\overline{\langle PILK^* ME \rangle} \xrightarrow{\tau} \langle PI(L+1)K^* M'E[x \leftarrow a]$
	[ArrayUpdate]
	$I(L) = x[e'] \leftarrow e a \stackrel{\text{def}}{=} E(x) M E e' \to n M E e \to v$
[Update]	$M(a) = [v_0, \ldots, v_m] 0 \le n \le m$
$I(L) = x \leftarrow e x \in dom(E) M E e \to v$	$M' \stackrel{\text{def}}{=} M[a \leftarrow [v_0, \dots, v_m] \{v_n / v\}]$
$\langle PILK^*ME\rangle \xrightarrow{\tau} \langle PI(L+1)K^*ME[x\leftarrow v]\rangle$	$\langle PILK^*ME\rangle \xrightarrow{\tau} \langle PI(L+1)K^*M'E\rangle$
[ReA	AD] [Print]
$I(L) = \operatorname{read} x$	$I(L) = \mathbf{print} \ e M \ E \ e \to lit$
$\langle PILK^*ME\rangle \xrightarrow{\text{read } lit} \langle PI(L+1)K^*ME[x \leftarrow lit)$	$[t]\rangle \qquad \langle PILK^*ME\rangle \xrightarrow{\text{print lit}} \langle PI(L+1)K^*ME\rangle$
[BranchT]	[BranchF]
$I(L) =$ branch $e L_1 L_2$ $M E e \rightarrow$ true	$I(L) =$ branch $e L_1 L_2$ $M E e \rightarrow$ false
$\langle PILK^* ME \rangle \xrightarrow{\tau} \langle PIL_1K^* ME \rangle$	$\langle PILK^* ME \rangle \xrightarrow{\tau} \langle PIL_2K^* ME \rangle$
[Gото]	[Stop]
I(L) = goto L'	I(L) = stop
$\langle PILK^* ME \rangle \xrightarrow{\tau} \langle PIL'K^* ME \rangle$	$\langle P I L K^* M E \rangle \xrightarrow{\text{stop}} \langle P \emptyset L K^* M E \rangle$
[CALL]	
$I(L) = \operatorname{call} x = e(e_1, \dots, e_n)$ M E e \rightarrow F	
$P(F) = F(x_1, \ldots, x_n) : D_F I' \stackrel{\text{def}}{=} P(F, \text{ active})$	[Return]
$L' \stackrel{\text{def}}{=} \operatorname{start}(I') M E [x_1 = e_1, \dots, x_n = e_n] \rightsquigarrow E'$	$I(L) =$ return $e M E e \rightarrow v$
$\overline{\langle PILK^*ME\rangle \xrightarrow{\tau} \langle PI'L'(K^*, \langle I(L+1)xE\rangle)ME'\rangle}$	$\overline{\langle PIL(K^*, \langle I'L'xE'\rangle) ME\rangle \stackrel{\tau}{\longrightarrow} \langle PI'L'K^*ME'[x\leftarrow v]\rangle}$
[AssumePass]	[AssumeDeopt]
$I(L) =$ assume e^* else $\xi \ \tilde{\xi}^* \forall m, M E \ e_m \rightarrow $ true	$I(L) =$ assume e^* else $\xi \ \tilde{\xi}^* \neg (\forall m, M E e_m \rightarrow \text{true})$
$\overline{\langle PILK^* ME \rangle} \xrightarrow{\tau} \langle PI(L+1)K^* ME \rangle$	$\overline{\langle PILK^* ME \rangle \stackrel{\tau}{\longrightarrow} deoptimize(\langle PILK^* ME \rangle, \xi, \tilde{\xi}^*)}$
	[DeoptimizeConf]
$M E V A \rightsquigarrow E' I' \stackrel{\text{def}}{=} P(F', V')$	
$\forall q \in 1,, r,$	
$\tilde{\xi}_q = F_q \cdot V_q \cdot L_q x_q V A_q$	
$M E VA_q \rightsquigarrow E_q I_q \stackrel{\text{def}}{=} P(F_q, V_q)$	$K_q \stackrel{\text{der}}{=} \langle I_q \ L_q \ x_q \ E_q \rangle$
deoptimize($\langle PILK^*ME \rangle$, $F'.V'.L'VA$,	$\tilde{\xi}_1, \ldots, \tilde{\xi}_r$) $\stackrel{\text{def}}{=} \langle P I' L' (K^*, K_1, \ldots, K_r) M E' \rangle$

[EVALENV] $\frac{M E \ e_1 \rightarrow v_1 \quad \dots \quad M E \ e_n \rightarrow v_n}{M E \ [x_1 = e_1, \ \dots, \ x_n = e_n] \rightsquigarrow [x_1 \rightarrow v_1, \ \dots, \ x_n \rightarrow v_n]}$

Fig. 17. Reduction relation $C \xrightarrow{\tau} C'$ for sourir IR.

Proceedings of the ACM on Programming Languages, Vol. 2, No. POPL, Article 49. Publication date: January 2018.

LEMMA 5.2 (TRANSITIVITY). If R_{12} is a weak bisimulation between P_1 and P_2 , and R_{23} is a weak bisimulation between P_2 and P_3 , then the composed relation $R_{13} \stackrel{\text{def}}{=} (R_{12}; R_{23})$ is a weak bisimulation between P_1 and P_3 .

Definition 5.3 (Version bisimilarity). Let V_1 , V_2 be two versions of a function F in P, and let $I_1 \stackrel{\text{def}}{=} P(F, V_1)$ and $I_2 \stackrel{\text{def}}{=} P(F, V_2)$. V_1 and V_2 are (weakly) bisimilar if $\langle P I_1 \operatorname{start}(I_1) K^* M E \rangle$ and $\langle P I_2 \operatorname{start}(I_2) K^* M E \rangle$ are weakly bisimilar for all K^* , M, E.

Definition 5.4 (Equivalence). P_1 , P_2 are equivalent if start(P_1), start(P_2) are weakly bisimilar.

5.4 Deoptimization Invariants

We can now give a formal definition of the invariants from Section 3.2: Version Equivalence holds if any pair of versions (V_1, V_2) of a function F are bisimilar; Assumption Transparency holds if for any configuration C, at an **assume** e^* else $\xi \tilde{\xi}^*$, C, is bisimilar to deoptimize $(C, \xi, \tilde{\xi}^*)$, as defined in Figure 17, DEOPTIMIZECONF.

5.5 Creating Fresh Versions and Injecting Assumptions

Configuration *C* is *over* location *F*.*V*.*L* if it is $\langle PP(F, V) LK^* ME \rangle$, where P(F, V) denotes the instructions at version *V* of *F* in *P*. Let $C[F.V.L \leftarrow F'.V'.L']$ be the configuration $\langle PP(F', V') L'K^* ME \rangle$. More generally, $C[X \leftarrow Y]$ replaces various components of *C*. For example, $C[P_1 \leftarrow P_2]$ updates the program in *C*; if only the versions change between two locations *F*.*V*.*L* and *F*.*V'*.*L*, write $C[V \leftarrow V']$ instead of repeating the locations, etc.

THEOREM 5.5. Creating a new copy of the currently active version of a function, possibly adding new assume instructions, returns an equivalent program.

PROOF. Consider P_1 with a function F with active version V_1 . Adding a version yields P_2 with new active version V_2 of F such that

- any label L of V_1 exists in V_2L : the instruction at L in V_1 and V_2 are identical except for assume instructions updated so that **assume** e^* else $\xi \ \tilde{\xi}^*$ in V_1 has a corresponding **assume** e^* else $F.V_1.L$ ld in V_2 where ld is the identity over the environment at L.
- V_2 may contain extra empty assume instructions: for any instruction *i* at *L* in V_1 , V_2 may contain an assume of the form **assume true else** $F.V_1.L$ ld, where ld is the identity mapping over the environment at *L*, followed by *i* at a fresh label *L'*.

Let us write I_1 and I_2 for the instructions of V_1 and V_2 respectively. Stack K_2^* is a *replacement* of K_1^* if it is obtained from K_1^* by replacing continuations of the form $\langle I_1 \ L \ x \ E \rangle$ by $\langle I_2 \ L \ x \ E \rangle$. Replacement is a device used in the proof and does not correspond to any of the reduction rules. We define a relation *R* as the smallest relation such that :

- (1) For any configuration C_1 over P_1 , R relates C_1 to $C_1[P_1 \leftarrow P_2]$.
- (2) For any configuration C_1 over a $F.V_1.L$ such that L in V_2 is not an added assume, R relates C_1 to $C_1[P_1 \leftarrow P_2][V_1 \leftarrow V_2]$.
- (3) For any configuration C_1 over a $F.V_1.L$ such that at L in V_2 is a newly added assume followed by label L', R relates C_1 to both (a) $C_1[F.V_1.L \leftarrow F.V_2.L]$ and (b) $C_1[F.V_1.L \leftarrow F.V_2.L']$.
- (4) For any related pair $(C_1, C_2) \in R$, where K_1^* is the call stack of C_2 , for any replacement K_2^* , the pair $(C_1, C_2[K_1^* \leftarrow K_2^*])$ is in R.

The proof proceeds by showing that *R* is a bisimulation. If a related pair $(C_1, C_2) \in R$ comes from the cases (1), (2) or (3) of the definition of *R*, we say that it is a *base pair*. A pair (C_1, C_2) in case (4) is defined from another pair $(C_1, C_2) \in R$, such that the call stack of C_2 is a replacement of the

stack of C'_2 . If $(C_1, C'_2) \in R$ is a base pair, we say that it is the base pair of (C_1, C_2) . Otherwise, we say that the base pair of (C_1, C_2) is the base pair of (C_1, C'_2) .

Bisimulation proof: generalities. To prove that R is a bisimulation, consider all related pairs $(C_1, C_2) \in R$ and show that a reduction from C_1 can be matched by C_2 and conversely. Without loss of generality, assume that C_2 is not a newly added assume instruction – that the base pair of (C_1, C_2) is not in the case (3,b) of the definition of R. Indeed, the proof of the case (3,b) follows from proof of the case (3,a). In the case (3,b), C_2 is a newly added assume instruction **assume true else** ... at L followed by L'. C_2 can only reduce silently into $C'_2 \stackrel{\text{def}}{=} C_2[L \leftarrow L']$, which is related to C_1 by the case (3,a). The empty reduction sequence from C_1 matches this reduction from C_2 . Conversely, assume the result in the case (3,a), then any reduction of C_1 can be matched from C'_2 , and thus matched from C_2 by prepending the silent reduction $C_2 \stackrel{\tau}{\longrightarrow} C'_2$ to the matching reduction sequence. Finally, if (C_1, C_2) comes from case (4) and has a base pair (C_1, C'_2) from (3,b), and C_2 has label L followed by L', then the bisimulation property for $(C_1, C_2) \in R$ comes from the one of $(C_1, C_2[L \leftarrow L']) \in R$ by the same reasoning.

Bisimulation proof: easy cases. The easy cases of the proof are the reductions $C_1 \xrightarrow{A_\tau} C'_1$ where neither C_1 nor C'_1 are over V_1 , and the reductions $C_2 \xrightarrow{A_\tau} C'_2$ where neither C_2 nor C'_2 are over V_2 . For $C_1 \xrightarrow{A_\tau} C'_1$, define C'_2 as $C'_1[P_1 \leftarrow P_2]$, and both $C_2 \xrightarrow{A_\tau} C'_2$ and $(C'_1, C'_2) \in R$ hold. The $C_2 \xrightarrow{A_\tau} C'_2$ case is symmetric, defining C'_1 as $C'_2[P_2 \leftarrow P_1]$.

Bisimulation proof: harder cases. The harder cases are split in two categories: version-change reductions (deoptimizations, functions call and returns), and same-version reductions within V_1 in P_1 or V_2 in P_2 . We consider same-version reductions first. Without loss of generality, assume that the pair $(C_1, C_2) \in R$ is a base pair, that is a pair related by the cases (2) or (3) of the definition of R, but not (4) – the case that changes the call stack of the configuration. Indeed, if pair $(C_1, C_2) \in R$ comes from (4), the only difference between this pair and its base pair $(C_1, C_2) \in R$ is in the call stack of C_2 and C'_2 . This means that C_2 and C'_2 have the exact same reduction behavior for non-version-change reductions. As long as the proof that the related configurations C_1 and C_2 match each other does not use version-change reductions (a property that holds for the proofs of the non-version-change reduction (deoptimization, call or return), prove that it can be matched from C_1 by reasoning on whether C_2 or C'_2 are assume instructions, coming from V_1 or newly added.

- If none of them are assume instructions, then they are both in the case (2) of the definition of R, they are equal to $C_1[V_1 \leftarrow V_2]$ and $C'_1[V_1 \leftarrow V_2]$ respectively, so $C_1 \xrightarrow{A_\tau} C'_1$ and $(C'_1, C'_2) \in R$ hold.
- If C_2 or C'_2 are assume instructions coming from V_1 , the same reasoning holds the problematic case where the assume is C_2 and the guards do not pass is not considered here as the reduction is not a deoptimization.
- If C'_2 is a newly added assume in V_2 at L followed by L', C_2 is an instruction of V_2 copied from V_1 , so (C_1, C_2) are in the case (2) of the definition of R and C_1 is $C_1[V_2 \leftarrow V_1]$. The reduction from C_2 corresponds to a reduction $C_1 \xrightarrow{A_\tau} C'_1$ in P_1 with $C'_1 \stackrel{\text{def}}{=} C'_2[V_2 \leftarrow V_1]$, and $(C'_1, C'_2) \in R$ by the case (3,a) of the definition of R.

The reasoning for transitions $C_1 \xrightarrow{A_r} C'_1$ that have to be matched from C_2 and are not version-change transitions (deoptimization, function calls or return) is similar. C_2 cannot be a new assume, so we

have $C_2 \xrightarrow{A_r} C'_2$, and either C'_2 is not a new assume and matches C_1 by case (2) of the definition of R, or it is a new assume and it matches it by the case (3,a).

Bisimulation proof: final cases. The cases that remain are the hard cases of version-change reductions: function call, return and deoptimization. If $C_1 \xrightarrow{A_{\tau}} C'_1$ is a deoptimization reduction, then C_1 is over a location $F.V_1.L$ in P_1 , and its instruction is **assume** e^* else $\xi \ \tilde{\xi}^*$, and C'_1 is deoptimize($C_1, \xi, \tilde{\xi}^*$). C_2 is over the copied instruction **assume** e^* else $F.V_1.L$ Id and Id is the identity. C_2 also deoptimizes, given that the tests give the same results in the same environment, so we have $C_2 \xrightarrow{\tau} C'_2$ for $C'_2 \stackrel{\text{def}}{=}$ deoptimize($C_2, F.V.L_1 \ \text{Id}, \emptyset$). C'_2 is over $F.V_1.L$, that is the same assume instruction as C_1 , so it also deoptimizes, to $C''_2 \stackrel{\text{def}}{=}$ deoptimize($C'_2, \xi, \tilde{\xi}^*$). We show that C'_1 and C''_2 are related by R:

- If $(C_1, C_2) \in R$ is a base pair, then C_1 is $C_2[V_2 \leftarrow V_1]$. In particular, the two configurations have the same environment, and C'_2 is identical to C_2 except it is over $F.V.L_1$. It is thus equal to C_1 . As a consequence, C'_1 and C''_2 , which are obtained from C_1 and C'_2 by the same deoptimization reduction, are the same configurations, and related in R.
- If C_1 and C_2 are related by the case (4) of the definition of R, the stack of C_2 is a replacement of the stack of C_1 . The same reasoning as in the previous case shows that configurations C'_1 and C''_2 are identical, except that the stack of C''_2 is a replacement of the stack of C'_1 : they are related by the case (4) of the definition of R.

Conversely, if $C_2 \xrightarrow{A_{\tau}} C'_2$ is a deoptimization instruction then, by the same reasoning as in the proof of matching a deoptimization of C_1 , C'_2 is identical to C_1 (modulo replaced stacks). This means that the empty reduction sequence from C_1 matches the reduction of C_2 .

If $C_1 \xrightarrow{A_{\tau}} C'_1$ is a function call transition,

$$\langle P_1 I_1 L K_1^* M E \rangle \xrightarrow{\tau} \langle P_1 I_1' L' (K^*, \langle I_1 (L+1) \times E \rangle) M E' \rangle$$

 C_2 is on the same call with the same arguments, so it takes a transition $C_2 \xrightarrow{\tau} C'_2$ of the form

$$\langle P_2 I_2 L K_2^* M E \rangle \xrightarrow{\tau} \langle P_2 I_2' L' (K^*, \langle I_2 (L+1) \mathbf{x} E \rangle) M E' \rangle$$

The stack of C'_2 is a replacement of the stack of C'_1 : assuming that K^*_2 is a replacement of K^*_1 , the difference in the new continuation is precisely the definition of stack replacement – note that it is precisely this reasoning step that required the addition of case (4) in the definition of R. Also, the new instruction streams I'_1 and I'_2 are either identical (if the function is not F itself) or equal to I_1 and I_2 respectively, so we do have $(C'_1, C'_2) \in R$ as expected. The proof of the symmetric case, matching a function call from C_2 , is identical.

If $C_1 \xrightarrow{A_\tau} C'_1$ is a function return transition

$$\langle P_1 I_1 L (K^*, \langle I'_1 L' x E' \rangle) M E \rangle \stackrel{\tau}{\longrightarrow} \langle P_1 I'_1 L' K^*_1 M E' [x \leftarrow \nu] \rangle$$

then $C_2 \xrightarrow{A_r} C'_2$ is also a function return transition

$$\langle P_2 I_2 L (K^*, \langle I'_2 L' x E' \rangle) M E \rangle \stackrel{\tau}{\longrightarrow} \langle P_2 I'_2 L' K^*_2 M E' [x \leftarrow \nu] \rangle$$

We have to show that C'_1 and C'_2 are related by R. The environments and heaps of the two configurations are identical. We know that the stack of C_2 is a replacement of the stack of C_1 , which means that K_2^* a replacement of K_1^* , and that either I'_1 and I'_2 are identical or they are respectively equal to I_1 and I_2 . In either case, C'_1 and C'_2 are related by R. The proof of the symmetric case, matching a function return from C_2 , is identical. We have established that R is a bisimulation. Finally, remark that our choice of *R* also proves that the new version respects the assumption transparency invariant. A new assume at *L* in V_2 is of the form **assume true else** *F*. V_1 .*L* ld, with ld the identity environment. Any configuration *C* over *F*. V_2 .*L* is related by R^{-1} to $C[F.V_2.L \leftarrow F.V_1.L]$, which is equal to deoptimize($C, F.V_1.L \operatorname{Id}, \emptyset$). These two configurations are related by the bisimulation R^{-1} , so they are bisimilar.

LEMMA 5.6. Adding a new predicate e' to an existing assume instruction **assume** e^* else $\xi \tilde{\xi}^*$ of P_1 returns an equivalent program P_2 .

PROOF. This is a consequence of the invariant of assumption transparency. Let R_{P_1} be the bisimilarity relation for configurations over P_1 , and F.V.L be the location of the modified assume. Let us define the relation R between P_1 and P_2 by

$$(C_1, C_2) \in R \quad \iff \quad (C_1, C_2[P_2 \leftarrow P_1]) \in R_{P_1}$$

We show that *R* is a bisimulation. Consider $(C_1, C_2) \in R$. If C_2 is not over *F.V.L*, the reductions of C_2 (in P_2) and $C_2[P_2 \leftarrow P_1]$ (in P_1) are identical, and the latter configuration is, by assumption, bisimilar to C_1 , so it is immediate that any reduction from C_1 can be matched by C_2 and conversely. If C_2 is over *F.V.L*, we can compare its reduction behavior (in P_2) with the one of $C_2[P_2 \leftarrow P_1]$ (in P_1). The first configuration deoptimizes when one of the e^* , e' is not true in the environment of C_2 , while the second deoptimizes when one of the e^* is not true – in the same environment. If C_2 gives the same boolean value to both series of test, then the two configurations have the same reduction behavior, and (C_1, C_2) match each other by the same reasoning as in the previous paragraph. The only interesting case is the configurations C_2 that pass all the tests in e^* , but fail e'. Let us show that, even in that case, the reductions of C_1 and C_2 match each other. The following diagram will be useful to follow the proof below:



Let us first show that the reductions of C_2 can be matched by C_1 . The only possible reduction from C_2 , given our assumptions, is $C_2 \xrightarrow{\tau} deoptimize(C_2, \xi, \tilde{\xi}^*)$. We claim that the empty reduction sequence from C_1 matches it, that is, that $(C_1, deoptimize(C_2, \xi, \tilde{\xi}^*)) \in R$. By definition of R, this goal means that C_1 and deoptimize $(C_2, \xi, \tilde{\xi}^*)[P_2 \leftarrow P_1]$ are bisimilar in P_1 . But the latter configuration is the same as deoptimize $(C_2[P_2 \leftarrow P_1], \xi, \tilde{\xi}^*)$, which is bisimilar to C_2 by the invariant of assumption transparency, and thus to C_1 by transitivity. Conversely, we show that the reductions of C_1 can be matched by C_2 . Suppose a reduction $C_1 \xrightarrow{A_\tau} C'_1$. The configuration deoptimize $(C_2, \xi, \tilde{\xi}^*)[P_2 \leftarrow P_1]$ is bisimilar to C_1 (same reasoning as in the previous paragraph), so there is a matching state C''_1 such that deoptimize $(C_2, \xi, \tilde{\xi}^*)[P_2 \leftarrow P_1] \xrightarrow{A_\tau} C''_1$ in P_1 with $(C'_1, C''_1) \in R_{P_1}$. We can transpose this reduction in P_2 : deoptimize $(C_2, \xi, \tilde{\xi}^*) \xrightarrow{A_\tau} C''_1[P_1 \leftarrow P_2]$ in P_2 , and thus $C_2 \xrightarrow{A_\tau} C''_1[P_1 \leftarrow P_2]$. This matches the reduction of C_1 , given that our assumption $(C'_1, C''_1) \in R_{P_1}$ exactly means that $(C'_1, C''_1[P_1 \leftarrow P_2]) \in R$.

Proceedings of the ACM on Programming Languages, Vol. 2, No. POPL, Article 49. Publication date: January 2018.

6 OPTIMIZATION CORRECTNESS

The proofs of the optimizations from Section 4 are easier than the proofs for deoptimization invariants in the previous section (although, as program transformations, they seem more elaborate). This comes from the fact that the classical optimizations rewrite an existing version and interact little with deoptimization.

6.1 Constant Propagation

We say that given a version V, a *static environment* SE for label L maps a subset of the variables in scope at L to values. A static environment is *valid*, written $SE \models L$, if for any configuration C over L reachable from the start of V we have that SE is a subset of the lexical environment E. Constant propagation can use a classic work-queue data-flow algorithm to compute a valid static environment SE at each label L. It then replaces, in the instruction at L, each expression or simple expression that can be evaluated in SE by its value. This is speculative since assumption predicates of the form x = lit populate the static environment with the binding $x \to lit$.

LEMMA 6.1. For any version V_1 , let V_2 be the result of constant propagation. V_1 and V_2 are bisimilar.

PROOF. The relation R to use here for bisimulation is the one that relates each reachable C_1 in reachable(P_1) to the corresponding state $C_2 \stackrel{\text{def}}{=} C_1[V_1 \leftarrow V_2]$ in reachable(P_2). Consider two related C_1, C_2 over L, and SE be the valid static environment at L inferred by our constant propagation algorithm. Reducing the next instruction of C_1 and C_2 will produce the same result, given that they only differ by substitutions of subexpressions by values that are valid under the static environment SE, and thus under E. If $C_1 \stackrel{A_{\tau}}{\longrightarrow} C'_1$ then $C_2 \stackrel{A_{\tau}}{\longrightarrow} C'_2$, and conversely.

The restriction of our bisimulation *R* to reachable configurations introduced is crucial for the proof to work. Indeed, a configuration that is not reachable may *not* respect the static environment *SE*. Consider the following example, with V1 on the left and V2 on the right.

L1	$\mathbf{var} \mathbf{x} = 1$	L1	var x = 1
	print x + x		print 2
	return 3		return 3

Now consider a pair of configurations at L1 with the binding $x \rightarrow 0$ in the environment.

$$C_1 \stackrel{\text{def}}{=} \langle P P(F, V_1) L_1 K^* M [\mathbf{x} \to 0] \rangle \qquad \qquad C_2 \stackrel{\text{def}}{=} \langle P P(F, V_2) L_1 K^* M [\mathbf{x} \to 0] \rangle$$

They would be related by the relation R used by the proof, yet they are not bisimilar: we have $C_1 \xrightarrow{\text{print 0}} C'_1$ as the only transition of C_1 in V_1 , and $C_2 \xrightarrow{\text{print 2}} C'_2$ as the only transition of C_2 in V_2 .

6.2 Unreachable Code Elimination

The following two lemmas are trivial: the simple version-change mapping between configurations on the two version is clearly a bisimulation. In the first case, this comes from the case that **branch true** L_1 L_2 and **goto** L_1 reduce in the example same way. In the second case, unreachable configurations are not even considered by the proof.

LEMMA 6.2. Replacing branch true $L_1 L_2$ by goto L_1 or branch false $L_1 L_2$ by goto L_2 results in an equivalent program.

LEMMA 6.3. Removing an unreachable label results in an equivalent program.

Proceedings of the ACM on Programming Languages, Vol. 2, No. POPL, Article 49. Publication date: January 2018.

6.3 Function Inlining

Assume that the function *F* has active version Vcallee. If the new version contains a call to *F*, **call** res = $F(e_1, ..., e_n)$ with return label Lret (the label after the call), inlining removes the call and instead:

- declares a fresh mutable return variable **var** res = **nil**;
- for the formal variables x_1 , ..., of F, defines the argument variables **var** $x_1 = se_1, ..., var x_n = se_n$;
- inserts the instructions from Vcallee, replacing each instruction **return** *e* by the sequence: res $\leftarrow e$; **drop** x_1 ; ...; **drop** x_n ; **goto** Lret

THEOREM 6.4. The inlining transformation presented returns a version equivalent to the caller version.

PROOF. The key idea of the proof is that any environment *E* in the inlined instruction stream can be split into two disjoint parts: an environment corresponding to the caller function, E_{caller} , and an environment corresponding to the callee, E_{callee} . To build the bisimulation, we relate the inlined version, on one hand, with the callee on the other hand, *when* the callee was called by the called at the inlined call point. This takes two forms:

- If a configuration is currently executing in the callee, and has the caller on the top of the call stack with the expected return address, we relate it to a configuration in the inlined version (at the same position in the callee). The environment of the inlined version is exactly the union of the callee environment (the environment of the configuration) and the caller environment (found on the call stack).
- If the stack contains a caller frame above a callee frame, we relate this to a single frame in the inlined version; again, there is a bidirectional correspondence between inlined environment and a pair of a caller and callee environment.

To check that this relation is a bisimulation, there are three interesting cases:

- If a transition is purely within the callee's code on one side, and within the inlined version of the callee on the other, it suffices to check that the environment decomposition is preserved. During the execution of inlinee, E_{caller} never changes, given that the instruction coming from the callee do not have the caller's variable in scope—and thus cannot mutate them.
- If the transition is a call of the callee from the caller on one side, and the entry into the declaration of the return variable **var** res = **nil** on the other, we step through the silent transitions that bind the call parameters **var** $x_1 = e_1, ..., var x_n = e_n$ and get to a state in the inlined function corresponding to the start of the callee.
- If the transition is a **return** *e* of the callee to the caller on one side, and the entry into the result assignment res ← *e* on the other, we similarly step through the **drop** *x* for each *x* in the callee's environment, and get to related state on the label *ret* following the function call.

6.4 Unrestricted Deoptimization

Consider P_1 containing an assume at L_1 , followed by i_m at $L_2 \stackrel{\text{def}}{=} (L_1 + 1)$. Let i_m be such it has a unique successor, is the unique predecessor of L_1 , and is not a function call, has no side-effect, does not modify the heap (array write or creation), and does not modify the variables mentioned in the assume. Under these conditions, we can move the assume immediately after the successor of i_m . Let us name P_2 the program modified in this way.

LEMMA 6.5. Given a program P_1 , and P_2 obtained by permuting an assume instruction L_1 after i_m at L_2 under the conditions above, P_1 and P_2 are bisimilar.

PROOF. The applicability restrictions are specific enough that we can reason precisely about the structure of reductions around the permuted instructions. Consider a configuration C_1 over the assume at L_1 in P_1 , and the corresponding configuration $C_2 \stackrel{\text{def}}{=} C_1[P_1 \leftarrow P_2][L_1 \leftarrow L_2]$ over L_2 in P_2 . Instruction i_m has a single successor, so there is only one possible reduction rule. Since i_m is not an I/O instruction, it must be a silent action. Hence there is a unique C'_2 such that $C_2 \stackrel{A_{\tau}}{\longrightarrow} C'_2$ holds, and furthermore A_{τ} is τ . Configurations C_1 and C'_2 are over the same assume. Let E_1 and E_2 be environments of C_1 and C'_2 respectively, and E' be their common sub-environment that contain only the variables mentioned in the assume (i_m does not modify its variables). If all tests in the assume instruction are true under E', then C_1 and C'_2 silently reduce to C'_1 and C''_2 . C'_1 is over i_m at L_2 , so it reduces $C'_1 \stackrel{\tau}{\longrightarrow} C''_1$; notice that C''_1 and C''_2 are over the labels (L_2 +1) in P_1 and (L_1 +1) in P_2 , which are equal. If not all tests of the assume are true under E', then both C_1 and C'_2 deoptimize. The deoptimized configurations are the same

- their function, version and label are the same: the assume's deoptimization target;
- they have the same call stack: it only depends on the call stack of *C*₁ and the interpretation of the assume's extra frames under *E*';
- they have the same heap, as we assumed that *i_m* does not modify the heap;
- they have the same deoptimized environment: it only depends on *E*'.

Let us call C_0 the configuration resulting from either deoptimization transitions.

We establish bisimilarity using definition a relation R and proving it is a bisimulation. The following diagrams are useful to follow the definition of R and the proofs.



We define *R* as the smallest relation such that:

- (1) For any C_1 and C_2 as above, C_1 and C'_1 are related to C_2 .
- (2) For any C_1 and C_2 as above such that C_1 passes the assume tests (does not deoptimize), both C'_2 and C''_2 are related to C''_1 .
- (3) For any C over P_1 that is over neither L_1 nor L_2 , C and $C[P_1 \leftarrow P_2]$ are related.

We now prove that *R* is a bisimulation. Any pair of configurations that are not over either L_1 or L_2 come from the case (3), so they are identical and it is immediate that they match each other. The interesting cases are for matching pairs of configurations over L_1 or L_2 .

In the case where no deoptimization happens, the reductions in P_2 are either $C_2 \xrightarrow{\tau} C'_2$, where both configurations are related to C_1 , or $C'_2 \xrightarrow{\tau} C''_2$ which is matched by $C_1 \xrightarrow{\tau} C'_1$. The reductions in P_1 are either $C_1 \xrightarrow{\tau} C'_1$, which is matched by $C_2 \xrightarrow{\tau} C'_2 \xrightarrow{\tau} C''_2$ and $C'_2 \xrightarrow{\tau} C''_2$, or $C'_1 \xrightarrow{\tau} C''_1$, which are both related to C''_2 .

In the case where a deoptimization happens, the only reduction in P_1 is $C_1 \xrightarrow{\tau} C_0$, which is matched by $C_2 \xrightarrow{\tau} C'_2 \xrightarrow{\tau} C_0$ and $C'_2 \xrightarrow{\tau} C_0$. The reductions in P_2 are $C_2 \xrightarrow{\tau} C'_2$, which are matched by the empty reduction on C_1 and $C'_2 \xrightarrow{\tau} C_0$ are matched by $C_1 \xrightarrow{\tau} C_0$.

Finally, we show preservation of the assumption transparency invariant. We have to establish the invariant for P_2 , assuming the invariant for P_2 . We have to show that C_0 and C'_2 are bisimilar. C_0 is bisimilar to C_1 (this is the transparency invariant on P_1), and C_1 and C'_2 are bisimilar because they are related by the bisimulation R.

6.5 Predicate Hoisting

Hoisting predicates takes a version V_1 , an expression e, and two labels L_1, L_2 , such that the instruction at L_1, L_2 are both assume instructions and e is a part of the predicate list at L_1 . The pass copies e from L_1 to L_2 , if all variables mentioned in e are in scope at L_2 . If, after this step the e can be constant folded to **true** at L_1 by the optimization from Section 4.1, then it is removed from L_1 , otherwise the whole version stays unchanged.

LEMMA 6.6. Let V_2 be the result of hoisting e from L_1 to L_2 in V_1 . V_1 and V_2 are bisimilar.

PROOF. Copying is bisimilar due to the assumption transparency invariant and to the fact that the constant-folded version is bisimilar due to Lemma 6.1.

6.6 Assume Composition

Let V_1 , V_2 , V_3 be three versions of a function F with instruction streams I_1 , I_2 , I_3 , and L_1 , L_2 , L_3 labels, such that $I_1(L_1) =$ **assume** e_1 else $F.V_2.L_2$ VA_1 and $I_2(L_2) =$ **assume** e_2 else $F.V_3.L_3$ VA_2 . The composition pass creates a new program P_2 from P_1 identical but the assume $P_2(F.V_1.L_1)$ is replaced by **assume** e_1 , e_2 else $F.V_3.L_3$ $VA_2 \circ VA_1$ where ($[x_1 = e_1, ..., x_n = e_n] \circ VA$) is defined as $[x_1 = e_1 \{ \frac{VA(y)}{u} \forall y \in VA \}, ..., x_n = e_n \{ \frac{VA(y)}{u} \forall y \in VA \}].$

LEMMA 6.7. Let P_2 be the result of composing assume instructions at L_1 and L_2 . P_1 and P_2 are bisimilar.

PROOF. For $C_1 \xrightarrow{\tau} C'_1$, $C_2 \xrightarrow{\tau} C'_2$ over L_1 in P_1 , P_2 , we distinguish four cases:

- (1) If e_1 and e_2 both hold, the assume does not deoptimize in P_1 and P_2 and they behave identically.
- (2) If e_1 and e_2 both fail, the original program deoptimizes twice; the modified P_2 only once. Assuming deoptimizing under the combined varmap $MEVA_2 \circ VA_1 \rightsquigarrow E''$ produces an environment equivalent to $MEVA_1 \rightsquigarrow E'$ and $ME'VA_2 \rightsquigarrow E''$ the final configuration is identical. Since the extra intermediate step is silent, both programs are bisimilar.
- (3) If e_1 fails and e_2 holds, we deoptimize to V_3 in P_2 , but to V_2 in P_1 . As shown in case (2) the deoptimized configuration C'_2 over L_3 is equivalent to a post-deoptimization configuration of C'_1 , which, due to assumption transparency is bisimilar to C'_1 itself.
- (4) If e_1 holds and e_2 fails, deoptimize to V_3 in P_2 but not in P_1 . Again C'_2 is equivalent to a post-deoptimization state, which is, transitively, bisimilar to C'_1 .

Since a well-formed assume has only unique names in the deoptimization metadata, it is simple to show the assumption in (2) with a substitution lemma. $\hfill \Box$

7 DISCUSSION

Our formalization raises new questions and makes apparent certain design choices. In this section, we present insights into the design space for JIT implementations.

The Cost of Assuming. Assumptions restrict optimizations. Variables needed for deoptimization must be kept alive. Consider Figure 18, where an assume is at the end of a loop. As y is not modified, it can be removed. There is enough information to reconstruct it if needed. On the other hand, x cannot be synthesized out of thin air because it is computed in another function. Additionally, assume restricts code motion in two cases. First, side-effecting code cannot be

Lloop branch
$$z \neq 0$$
 Lbody Ldone
Lbody call $x = dostuff()$
var $y = x + 13$
assume e else $F.V.L$ [$x = x, y = x + 13$]
drop y
goto Lloop
Ldone ...

Fig. 18. Deoptimization keeps variables alive.

moved over an assume. Second, assume instructions cannot be hoisted over instructions that interfere with variables mentioned in metadata. It is possible to move assume forward, since data dependencies can be resolved by taking a snapshot of the environment at the original location. For the reverse effect, we support hoisting the predicate from one assume to another (see Section 4.5). Moving assume instructions up is tricky and also unnecessary, since in combination those two primitives allow moving checks to any position. In the above example, if *e* is invariant in the loop

a whole. *Lazy Deoptimization.* The runtime cost of an assume is the cost of monitoring the predicates. Suppose we speculate that the contents of an array remain unchanged throughout a loop. An implementation would have to check every single element of the array. An eager strategy where predicates are checked at every iteration is wasteful. It is more efficient to associate checks to operations that may invalidate the predicates, such as array writes, to invalidate the assumption, a strategy sometimes known as *lazy deoptimization*. We could implement dependencies by separating

body and there is an assume before Lloop, the predicate can be hoisted out of the loop. If the assume is only relevant for a subset of the instructions after the current location, it can be moved down as

assumptions from runtime checks. Specifically, let GUARDS[13] = true be the runtime check, where the global array GUARDS is a collection of all remote assumptions that can be invalidated by an operation, such as an array assignment. In terms of correctness, both eager and lazy deoptimization are similar; however, we would need to prove correctness of the dependency mechanism that modifies the global array.

Jumping Into Optimized Code. We have shown how to transfer control out of optimized code. The inverse transition, jumping into optimized code, is interesting as well. Consider executing the long running loop of Figure 19. The value of debug is constant in the loop, yet execution is stuck in the long running function and must branch on each iteration. A JIT can compile an optimized version that speculates on debug, but it may only use it on the next invocation. Ideally, the JIT would jump into the newly optimized code from the slow loop; this is known as hot loop transfer. Specifically, the next time Lo is reached, control is transferred to an equivalent location in the optimized version. To do so, continuation-passing style can be used to compile a staged continuation function from the beginning of the loop where debug is known to be false. The optimized continuation might look like cont in Figure 20. In some sense, this is easier than deoptimization because it strengthens assumptions rather than weakening them and all the values needed to construct the state at the target version are readily available.

stuc	k()
١	VЬ	ase

1	buse	
		call debug = debug()
	Lh	branch x < 1000000 Lo Lrt
	Lo	branch debug Lslow Lfast
	Lslow	
	Lfast	
		goto Lh
	Lrt	

Fig. 19. Long running execution.

cont(x)

```
Vopt
| Lh branch x < 1000000 Lfast Lrt
| Lfast ...
goto Lh
| Lrt ...
```



undo()	
Vs123	
Lo	assume e_1, e_2, e_3 else undo.Vs12.L0 []
Vs12	
Lo	assume e_1, e_2 else undo.Vs1.L0 []
Vs1	
Lo	assume e_1 else undo.Vbase.L0 []

Fig. 21. Undoing an isolated predicate.

Fine-Grained Deoptimization. Instead of blindly removing all assumptions on deoptimization, it is possible to undo only failing assumptions while preserving the rest. As shown in Figure 21, if e_2 fails in version Vs123, one can jump to the last version that did not rely on this predicate. By deoptimizing to version Vs1, assumption e_3 must be discarded. However, e_1 , e_3 still hold, so we would like to preserve optimizations based on those assumptions. Using the technique mentioned above, execution can be transferred to a version Ls13 that reintroduces e_3 . The overall effect is that we remove only the invalidated assumption and its optimizations. We are not aware of an existing implementation that explores such a strategy.

Simulating a Tracing JIT. A tracing JIT [Bala, Duesterwald, and Banerjia 2000; Gal, Eich, Shaver, Anderson, Mandelin, Haghighat, Kaplan, Hoare, Zbarsky, Orendorff, Ruderman, Smith, Reitmaier, Bebenita, Chang, and Franz 2009] records instructions that are executed in a trace. Branches and redundant checks can be discarded from the trace. Typically, a trace corresponds to a path through a hot loop. On subsequent



runs the trace is executed directly. The JIT ensures that execution follows the same path, otherwise it deoptimizes back to the original program. In this context Guo and Palsberg [2011] develop a framework for reasoning about optimizations applied to traces. One of their results is that dead store elimination is unsound, because the trace is only a partial view of the entire program. For example, a variable x might be assigned to within a trace, but never used. However, it is unsound to remove the assignment, because x might be used outside the trace. We can simulate their tracing formalism in sourir. Consider a variant of their running example shown in Figure 22, a trace of the loop **while** e ($x \leftarrow 0$; ...) embedded in a larger context. Instead of a JIT that records instructions, assume only branch targets are recorded. For this example, suppose the two targets Lbody and Ldone are recorded, which means the loop body executed once and then exited. In other words, the loop condition e was **true** the first time and **false** the second time. The compiler could unroll the loop twice and assert e for the first iteration and $\neg e$ for the second iteration (left). Then unreachable code elimination yields the right hand side, resembling a trace.

```
assume e else F.Vbase.Lloop [x = x, ...]
             branch e Lbody<sub>0</sub> Ldone
                                                                                      . . .
\mathsf{Lbody}_0
             \mathbf{x} \leftarrow \mathbf{0}
                                                                                     assume e else F.Vbase.Lloop [x = x, ...]
                                                                                     \mathbf{x} \leftarrow \mathbf{0}
             assume \neg e else F.Vbase.Lloop [x = x, ...]
                                                                                     . . .
             branch e Lbody1 Ldone
                                                                                     assume \neg e else F.Vbase.Lloop [x = x, ...]
Lbody<sub>1</sub>
             \mathbf{x} \leftarrow \mathbf{0}
                                                                                      . . .
             . . .
             goto Lloop
Ldone
```

Say x is not accessed after the store in this optimized version. In sourir, it is obvious why dead store elimination of x would be unsound: the deoptimization metadata indicates that x is needed for deoptimization and the store operation can only be removed it can be replayed. In this specific example, a constant propagation pass could update the metadata to materialize the write of 0, only when deoptimizing at the second assume. But, before the code can be reduced, loop unrolling might result in intermediate versions that are much larger than the original program. In contrast, tracing JITs can handle this case without the drastic expansion in code size [Gal et al. 2009], but lose more information about instructions outside of the trace.

8 CONCLUSIONS

Speculative optimizations are key to just-in-time optimization of dynamic languages. As these optimizations depend on predicates about the program state, the language implementation must monitor the validity of predicates and be ready to deoptimize the program if a predicate is invalidated. While, many modern compiler rely on this approach, the interplay between optimization and deoptimization often remains opaque.

Our contribution is to show that when the predicates and the deoptimization metadata are reified in the program representation, it becomes quite easy to define correct program transformations that are deoptimization aware. In this work we extend the intermediate representation with one new instruction, assume, which plays the double role of checking for the validity of predicates and specifying the actions required to deoptimize the program. Program transformations can inspect both the predicates that are being monitored and the deoptimization metadata and transform them when needed. The formalization presented here is for one particular intermediate language that we hope to be representative of a typical dynamic language. We present a bisimulation proof between multiple versions of the same function, optimized under different assumptions. We formalize deoptimization invariants between versions and show that they enable very simple proofs for standard compiler optimizations, constant folding, unreachable code elimination, and function inlining. We also prove correct three optimizations that are specifically dealing with deoptimizations, namely unrestricted deoptimization, predicate hoisting, and assume composition.

There are multiple avenues of future investigation. The optimizations presented here rely on intraprocedural analysis and the granularity of deoptimization is a whole function. If we were to extend this work to interprocedural analysis, it would become much trickier to determine what functions are to be invalidated as a speculation in one function may allow optimizations in many other functions. The current representation forces to check predicates before each use, but some predicates are cheaper to check by monitoring operations that could invalidate them. To do this would require changes to our model as the assume instruction would need to be split between a monitor and a deoptimization point. Lastly, the expressive power of predicates is an interesting question as there is a clear trade-off — richer predicates may allow more optimizations but are likely to be costlier to monitor.

ACKNOWLEDGMENTS

Jean-Marie Madiot provided guidance on the use and limitations of various notions of bisimulation. In particular, he suggested adding a non-silent stop transition to recover equi-termination from weak bisimilarity. Francesco Zappa Nardelli helped with the motivation and presentation of our work. We thank Sewell, Nardelli, Owens, Peskine, Ridge, Sarkar, and Strniša [2007] for writing and maintaining Ott. Our work supported by the National Science Foundation under Grants CCF–1544542, CCF–1318227, CCF–1618732, ONR award 503353, and the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation program (grant agreement 695412). Any opinions, findings, and conclusions expressed in this material may be those of the authors and likely do not reflect the views of our funding agencies.

REFERENCES

- Vasanth Bala, Evelyn Duesterwald, and Sanjeev Banerjia. 2000. Dynamo: A Transparent Dynamic Optimization System. In Programming Language Design and Implementation (PLDI). https://doi.org/10.1145/349299.349303
- Michael Bebenita, Florian Brandner, Manuel Fahndrich, Francesco Logozzo, Wolfram Schulte, Nikolai Tillmann, and Herman Venter. 2010. SPUR: A Trace-based JIT Compiler for CIL. In Conference on Object Oriented Programming Systems Languages and Applications (OOPSLA). https://doi.org/10.1145/1869459.1869517
- Clément Béra, Eliot Miranda, Marcus Denker, and Stéphane Ducasse. 2016. Practical validation of bytecode to bytecode JIT compiler dynamic deoptimization. *Journal of Object Technology (JOT)* 15, 2 (2016). https://doi.org/10.5381/jot.2016.15.2.a1 Project Chromium. 2017. V8 JavaScript Engine. https://chromium.googlesource.com/v8/v8.git.
- Daniele Cono D'Elia and Camil Demetrescu. 2016. Flexible on-stack replacement in LLVM. In *Code Generation and Optimization (CGO)*. https://doi.org/10.1145/2854038.2854061
- Stefano Dissegna, Francesco Logozzo, and Francesco Ranzato. 2014. Tracing Compilation by Abstract Interpretation. In *Principles of Programming Languages (POPL)*. https://doi.org/10.1145/2535838.2535866
- Gilles Duboscq, Thomas Würthinger, and Hanspeter Mössenböck. 2014. Speculation without regret: reducing deoptimization meta-data in the Graal compiler. In *Principles and Practices of Programming on the Java Platform (PPPJ)*. https://doi.org/ 10.1145/2647508.2647521
- Gilles Duboscq, Thomas Würthinger, Lukas Stadler, Christian Wimmer, Doug Simon, and Hanspeter Mössenböck. 2013. An Intermediate Representation for Speculative Optimizations in a Dynamic Compiler. In *Virtual Machines and Intermediate Languages (VMIL)*. https://doi.org/10.1145/2542142.2542143
- Stephen J. Fink and Feng Qian. 2003. Design, Implementation and Evaluation of Adaptive Recompilation with On-stack Replacement. In *Code Generation and Optimization (CGO)*. https://doi.org/10.1109/CGO.2003.1191549
- Andreas Gal, Brendan Eich, Mike Shaver, David Anderson, David Mandelin, Mohammad R. Haghighat, Blake Kaplan, Graydon Hoare, Boris Zbarsky, Jason Orendorff, Jesse Ruderman, Edwin W. Smith, Rick Reitmaier, Michael Bebenita, Mason Chang, and Michael Franz. 2009. Trace-based Just-in-time Type Specialization for Dynamic Languages. In Programming Language Design and Implementation (PLDI). https://doi.org/10.1145/1542476.1542528
- Shu-yu Guo and Jens Palsberg. 2011. The Essence of Compiling with Traces. In *Principles of Programming Languages (POPL)*. https://doi.org/10.1145/1926385.1926450
- Urs Hölzle, Craig Chambers, and David Ungar. 1992. Debugging Optimized Code with Dynamic Deoptimization. In Programming Language Design and Implementation (PLDI). https://doi.org/10.1145/143095.143114
- Kazuaki Ishizaki, Motohiro Kawahito, Toshiaki Yasue, Hideaki Komatsu, and Toshio Nakatani. 2000. A Study of Devirtualization Techniques for a Java Just-In-Time Compiler. In *Object-oriented Programming Systems, Language, and Applications* (OOPSLA). https://doi.org/10.1145/353171.353191
- Gary A. Kildall. 1973. A Unified Approach to Global Program Optimization. In *Principles of Programming Languages (POPL)*. https://doi.org/10.1145/512927.512945
- Xavier Leroy and Sandrine Blazy. 2008. Formal verification of a C-like memory model and its uses for verifying program transformations. *Journal of Automated Reasoning* 41, 1 (2008). https://doi.org/10.1007/s10817-008-9099-0
- Magnus O. Myreen. 2010. Verified Just-in-time Compiler on x86. In *Principles of Programming Languages (POPL)*. https://doi.org/10.1145/1706299.1706313
- Rei Odaira and Kei Hiraki. 2005. Sentinel PRE: Hoisting Beyond Exception Dependency with Dynamic Deoptimization. In *Code Generation and Optimization (CGO)*. https://doi.org/10.1109/CGO.2005.32
- Michael Paleczny, Christopher Vick, and Cliff Click. 2001. The Java Hotspot Server Compiler. In Java Virtual Machine Research and Technology (JVM). http://www.usenix.org/events/jvm01/full_papers/paleczny/paleczny.pdf
- Amr Sabry and Matthias Felleisen. 1992. Reasoning About Programs in Continuation-passing Style. In LISP and Functional Programming (LFP). https://doi.org/10.1145/141471.141563
- David Schneider and Carl Friedrich Bolz. 2012. The efficient handling of guards in the design of RPython's tracing JIT. In *Workshop on Virtual Machines and Intermediate Languages (VMIL)*. https://doi.org/10.1145/2414740.2414743
- Peter Sewell, Francesco Zappa Nardelli, Scott Owens, Gilles Peskine, Thomas Ridge, Susmit Sarkar, and Rok Strniša. 2007. Ott: Effective Tool Support for the Working Semanticist. In International Conference on Functional Programming (ICFP). https://doi.org/10.1145/1291151.1291155
- Sunil Soman and Chandra Krintz. 2006. Efficient and General On-Stack Replacement for Aggressive Program Specialization. In Software Engineering Research and Practice (SERP).
- Kunshan Wang, Yi Lin, Stephen M. Blackburn, Michael Norrish, and Antony L. Hosking. 2015. Draining the Swamp: Micro Virtual Machines as Solid Foundation for Language Development. In Summit on Advances in Programming Languages (SNAPL), Vol. 32. https://doi.org/10.4230/LIPIcs.SNAPL.2015.321
- Yudi Zheng, Lubomír Bulej, and Walter Binder. 2017. An Empirical Study on Deoptimization in the Graal Compiler. In European Conference on Object-Oriented Programming (ECOOP). https://doi.org/10.4230/LIPIcs.ECOOP.2017.30