

# Speculative Optimizations without Fear

Olivier Flückiger, Gabriel Scherer, Ming-Ho Yee, Aviral Goel,  
Jan Vitek, Amal Ahmed

Northeastern University, Boston, USA

September 12, 2017



1 Context

2 Sourir

3 Formalization

4 Optimizations

# Section 1

## Context

## Our work

Just-in-time (JIT) compilation is essential to efficient dynamic language implementations.

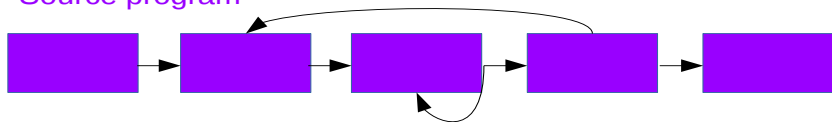
(Javascript, Lua, R... Java)

There is a blind spot in our formal understanding of JITs: speculation.

We present a language design to study speculative optimizations and prove them correct.

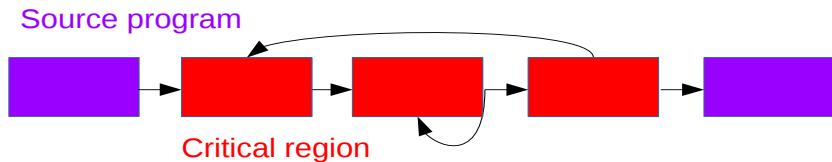
# Just-in-time compilation

Source program



JITs:

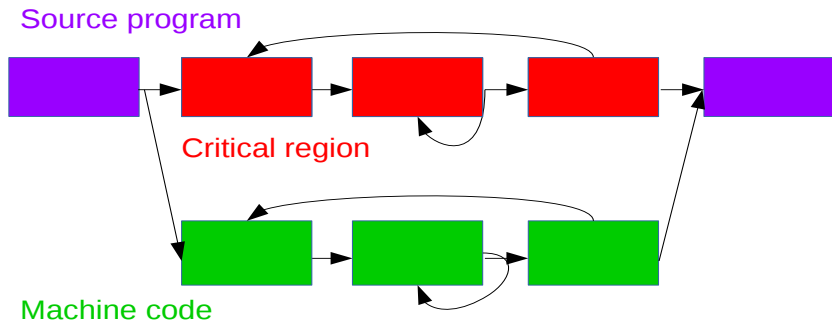
# Just-in-time compilation



JITs:

Profiling

# Just-in-time compilation

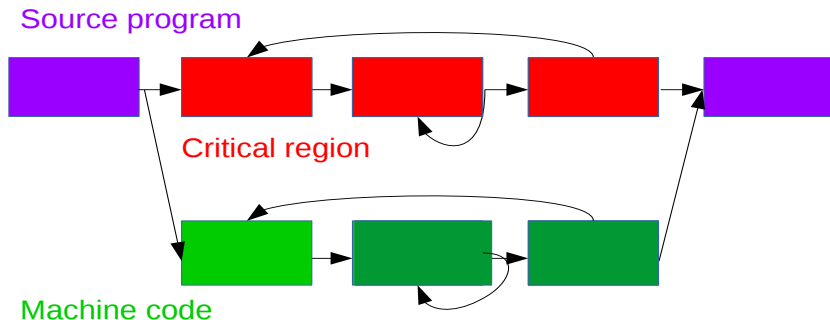


JITs:

Profiling

- + High/Low languages
- + Dynamic code generation/mutation

# Just-in-time compilation



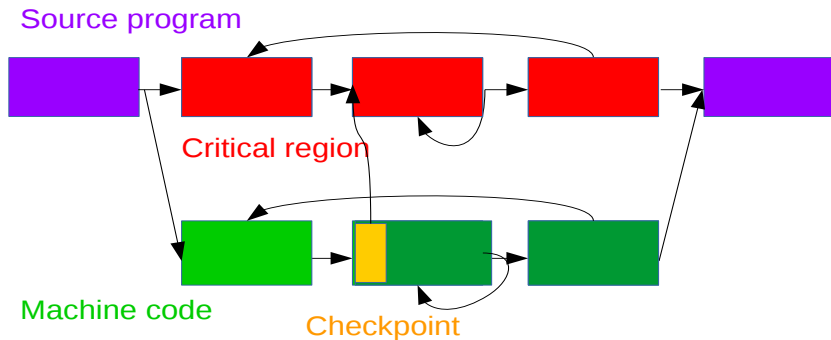
JITs:

Profiling

- + High/Low languages
- + Dynamic code generation/mutation
- + Speculation



# Just-in-time compilation



JITs:

Profiling

- + High/Low languages
- + Dynamic code generation/mutation
- + Speculation and bailout

## JITs:

- profiling
- high- and low-level languages (or multi-tiers, etc.)
- dynamic code generation + mutation
  
- speculation and bailout

## JITs:

- profiling
- high- and low-level languages (or multi-tiers, etc.)
- dynamic code generation + mutation  
eliminates interpretation overhead (constant factor)
- speculation and bailout

## JITs:

- profiling
- high- and low-level languages (or multi-tiers, etc.)
- dynamic code generation + mutation  
eliminates interpretation overhead (constant factor)
- speculation and bailout  
eliminates **dynamic features** overhead:  
dispatch (OO languages), type checks (Java),  
code loading (Java), redefinable primitives (R...)

## JITs:

- profiling
- high- and low-level languages (or multi-tiers, etc.)
- dynamic code generation + mutation  
eliminates interpretation overhead (constant factor)
- speculation and bailout  
eliminates **dynamic features** overhead:  
dispatch (OO languages), type checks (Java),  
code loading (Java), redefinable primitives (R...)

## JIT formalization: Myreen [2010]

- Stack language and x86 assembly
- dynamic code generation
- code mutation

mechanized in HOL!

## JITs:

- profiling
- high- and low-level languages (or multi-tiers, etc.)
- dynamic code generation + mutation  
eliminates interpretation overhead (constant factor)
- speculation and bailout  
eliminates **dynamic features** overhead:  
dispatch (OO languages), type checks (Java),  
code loading (Java), redefinable primitives (R...)

JIT formalization: [Myreen \[2010\]](#)

- Stack language and x86 assembly
- dynamic code generation
- code mutation

mechanized in HOL!

## JITs:

- profiling
- high- and low-level languages (or multi-tiers, etc.)
- dynamic code generation + mutation  
eliminates interpretation overhead (constant factor)
- speculation and bailout  
eliminates **dynamic features** overhead:  
dispatch (OO languages), type checks (Java),  
code loading (Java), redefinable primitives (R...)

JIT formalization: [Myreen \[2010\]](#)

- Stack language and x86 assembly
- dynamic code generation
- code mutation

mechanized in HOL!

What about **speculation**?

## JITs:

- profiling
- high- and low-level languages (or multi-tiers, etc.)
- dynamic code generation + mutation  
eliminates interpretation overhead (constant factor)
- speculation and bailout  
eliminates **dynamic features** overhead:  
dispatch (OO languages), type checks (Java),  
code loading (Java), redefinable primitives (R...)

JIT formalization: [Myreen \[2010\]](#)

- Stack language and x86 assembly
- dynamic code generation
- code mutation

mechanized in HOL!

What about **speculation**? [This work.](#)



## What do we want to know?

Speculation requires keeping **bailout data**.

How should optimizations maintain/transform bailout data?  
(inlining is tricky)

Does the presence of checkpoint restrict optimizations?  
(hoisting writes or IO is tricky)

When an assumption fails, how much of the other optimizations can keep?  
(non-stack-order is tricky)

How should practitioners reason about correctness?



DEEP SPECULATION.

THE OIL-SPECULATOR'S DREAM.

© 2001 HARP WEEK®

# Sourir

- high- and low-level languages
- dynamic code generation
- speculative optimization and bailout

# Sourir

- ~~high and low level languages~~
- dynamic code generation
- speculative optimization and bailout

a single bytecode language

# Sourir

- ~~high and low level languages~~ a single bytecode language
- ~~dynamic code generation~~ one unrolled multi-version program
- speculative optimization and bailout

# Sourir

- ~~high and low level languages~~ a single bytecode language
- ~~dynamic code generation~~ one unrolled multi-version program
- speculative optimization and bailout a **checkpoint** instruction

# Sourir

- ~~high and low level languages~~ a single bytecode language
- ~~dynamic code generation~~ one unrolled multi-version program
- speculative optimization and bailout a **checkpoint** instruction

$F_{fun}(c) \rightarrow$

$V_{tough} \rightarrow$

$L_0$  : **var**  $o = 1$

$L_1$  : **print** ( $c + o$ )

$V_{luck} \rightarrow$

$L_0$  : **assume** [( $c = 41$ )] **else**  $\langle F_{fun} \cdot V_{tough} \cdot L_1 [c = c, o = 1] \rangle$

$L_1$  : **print** 42

## Contribution

A language design to model speculative optimization: `Sourir`

A kit of correct program transformations and optimizations

A methodology to reason about correct speculative optimizations



## Section 2

Sourir

# A simple bytecode language

$i ::=$		$e ::=$
	<b>var</b> $x = e$	$se$
	<b>drop</b> $x$	$x[se]$
	$x \leftarrow e$	<b>length</b> ( $se$ )
	<b>array</b> $x[e]$	$primop(se^*)$
	<b>array</b> $x = [e^*]$	
	$x[e_1] \leftarrow e_2$	$se ::=$
	<b>branch</b> $e L_1 L_2$	$lit$
	<b>goto</b> $L$	$'F'$
	<b>print</b> $e$	$x$
	<b>read</b> $x$	
	<b>call</b> $x = e(e^*)$	$lit ::=$
	<b>return</b> $e$	$\dots, -1, 0, 1, \dots$
	<b>assume</b> $[e^*]$ <b>else</b> $\xi \tilde{\xi}^*$	<b>nil</b>   <b>true</b>   <b>false</b>
	<b>stop</b>	

# Versions

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

# Checkpoints

Checkpoint: **guards** + **bailout data**.

**assume**  $[(c = 41)]$  **else**  $\langle F_{fun} \cdot V_{tough} \cdot L_1 [c = c, o = 1] \rangle$

Guards: just a list of expressions returning booleans.

Bailout data:

- where to go:  $F_f \cdot V_w \cdot L_l$
- in what state:  $[x_1 = e_1, \dots, x_n = e_n]$
- (plus more: see inlining)

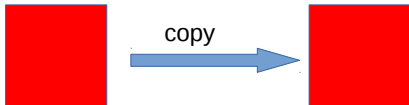
# Speculative optimization pipeline

Critical version

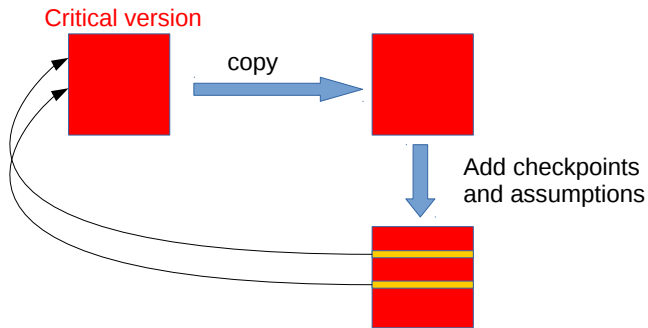


# Speculative optimization pipeline

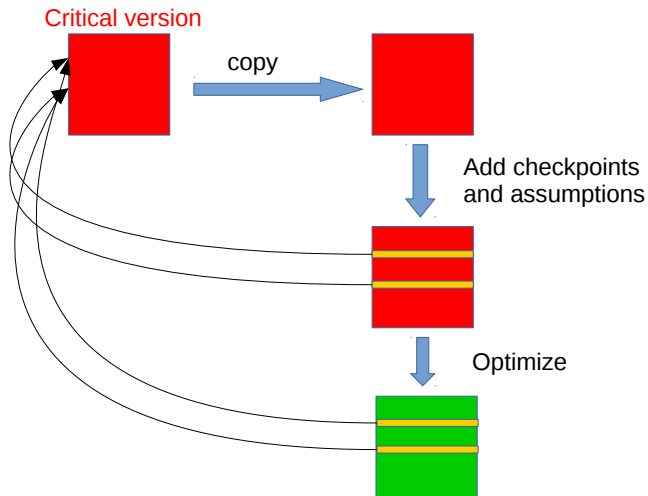
Critical version



# Speculative optimization pipeline

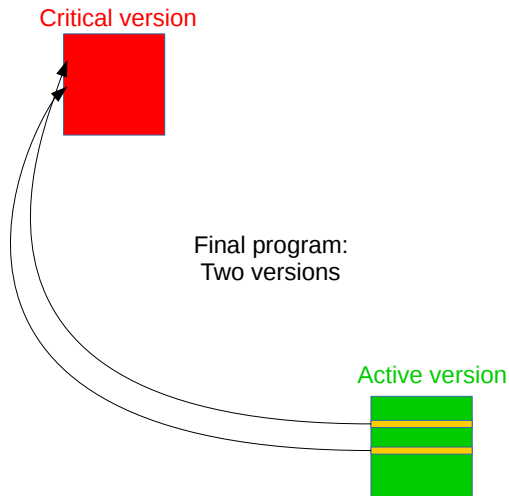


# Speculative optimization pipeline





# Speculative optimization pipeline



# Speculative optimization pipeline

Critical version



## Section 3

### Formalization

# Execution: Operational semantics

Configurations:

$$C ::= \langle P I L K^* M E \rangle$$

Actions:

$$A ::= \mathbf{read} \textit{ lit} \mid \mathbf{print} \textit{ lit}$$

$$A_\tau := A \mid \tau$$

$$T ::= A^*.$$

Reduction:

$$C_1 \xrightarrow{A_\tau^*} C_2$$

$$C_1 \xrightarrow{T} C_2$$



## Bailout invariants

Version invariant: All versions of a function are equivalent.  
(Necessary to replace the active version)

Bailout invariant: Bailing out **more** than necessary is correct.  
(Necessary to add new assumptions)

## Section 4

# Optimizations

## Branch pruning – from the kit

$V_{base} \rightarrow$

$L_1$  : **branch** ( $tag = INT$ )  $L_{int}$   $L_{nonint}$   
 $L_{int}$  : ...  
 $L_{nonint}$  : ...



## Branch pruning – from the kit

$V_{base} \rightarrow$

$L_1$  : **branch** ( $tag = INT$ )  $L_{int}$   $L_{nonint}$   
 $L_{int}$  : ...  
 $L_{nonint}$  : ...

$V_{opt} \rightarrow$

$L_0$  : **assume** [ $tag = INT$ ] **else**  $\langle F.V_{base}.L_1 \delta \rangle$   
 $L_1$  : **branch** ( $tag = INT$ )  $L_{int}$   $L_{nonint}$   
 $L_{int}$  : ...  
 $L_{nonint}$  : ...

Checkpoint + guard inserted

Bailout invariant!

## Branch pruning – from the kit

$V_{base} \rightarrow$

$L_1$  : **branch** ( $tag = INT$ )  $L_{int}$   $L_{nonint}$   
 $L_{int}$  : ...  
 $L_{nonint}$  : ...

$V_{opt} \rightarrow$

$L_0$  : **assume** [ $tag = INT$ ] **else**  $\langle F.V_{base}.L_1 \delta \rangle$   
 $L_1$  : **branch true**  $L_{int}$   $L_{nonint}$   
 $L_{int}$  : ...  
 $L_{nonint}$  : ...

constant folding

## Branch pruning – from the kit

$V_{base} \rightarrow$

$L_1$  : **branch** ( $tag = INT$ )  $L_{int}$   $L_{nonint}$   
 $L_{int}$  : ...  
 $L_{nonint}$  : ...

$V_{opt} \rightarrow$        $L_0$  : **assume** [ $(tag = INT)$ ] **else**  $\langle F.V_{base}.L_1 \delta \rangle$   
                   $L_{int}$  : ...

unreachable code elimination

# Inlining

$F_{main}() \rightarrow$

$V_{inlined} \rightarrow$

$L_0$  : **array**  $vec = [1, 2, 3, 4]$

$L_2$  : **var**  $size = nil$

$L_3$  : **var**  $obj = vec$

$L_{cp1}$  : **assume**  $[(obj \neq nil)]$  **else** ...

$L_5$  : **var**  $len = length(obj)$

$L_6$  :  $size \leftarrow (len * 4)$

$L_7$  : **drop**  $len$

$L_8$  : **drop**  $obj$

$L_9$  : **goto**  $L_{ret}$

$L_{ret}$  : **print**  $size$

$V_{base} \rightarrow \dots$

$F_{main}() \rightarrow$

$V_{base} \rightarrow$

$L_0$  : **array**  $vec = [1, 2, 3, 4]$

$L_2$  : **call**  $size = 'F_{size}(vec)$

$L_{ret}$  : **print**  $size$

,  $F_{size}(obj) \rightarrow$

$V_{opt} \rightarrow$

$L_{cp1}$  : **assume**  $[(obj \neq nil)]$  **else** ..

$L_{vec}$  : **var**  $len = length(obj)$

$L_3$  : **return**  $(len * 4)$

$V_{base} \rightarrow \dots$

# Inlining

$F_{main}() \rightarrow$

$V_{inlined} \rightarrow$

$L_0$  : **array**  $vec = [1, 2, 3, 4]$

$L_2$  : **var**  $size = nil$

$L_3$  : **var**  $obj = vec$

$L_{cp1}$  : **assume**  $[(obj \neq nil)]$  **else** ...

$L_5$  : **var**  $len = \mathbf{length}(obj)$

$L_6$  :  $size \leftarrow (len * 4)$

$L_7$  : **drop**  $len$

$L_8$  : **drop**  $obj$

$L_9$  : **goto**  $L_{ret}$

$L_{ret}$  : **print**  $size$

$V_{base} \rightarrow \dots$

$F_{main}() \rightarrow$

$V_{base} \rightarrow$

$L_0$  : **array**  $vec = [1, 2, 3, 4]$

$L_2$  : **call**  $size = 'F_{size}(vec)$

$L_{ret}$  : **print**  $size$

,  $F_{size}(obj) \rightarrow$

$V_{opt} \rightarrow$

$L_{cp1}$  : **assume**  $[(obj \neq nil)]$  **else** ..

$L_{vec}$  : **var**  $len = \mathbf{length}(obj)$

$L_3$  : **return**  $(len * 4)$

$V_{base} \rightarrow \dots$

Need for an extra frame in the inlined version:

$\langle F_{main} \cdot V_{base} \cdot L_{ret} \text{ size } [vec = vec] \rangle$

## Conclusion

All you need for speculation: versions + checkpoints.

Future work: bidirectional transformations.

Thanks!  
Questions?

Magnus O. Myreen. Verified just-in-time compiler on x86. In **Proceedings of the 37th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages**, POPL '10, pages 107–118, New York, NY, USA, 2010. ACM. ISBN 978-1-60558-479-9. doi: 10.1145/1706299.1706313.