

# Beyond ahead-of-time imperative compilation

Laure Gonnord & Matthieu Moy & Gabriel Radanne & other  
<https://compil-lyon.gitlabpages.inria.fr/>

Master 1, ENS de Lyon et Dpt Info, Lyon1

November 30, 2022

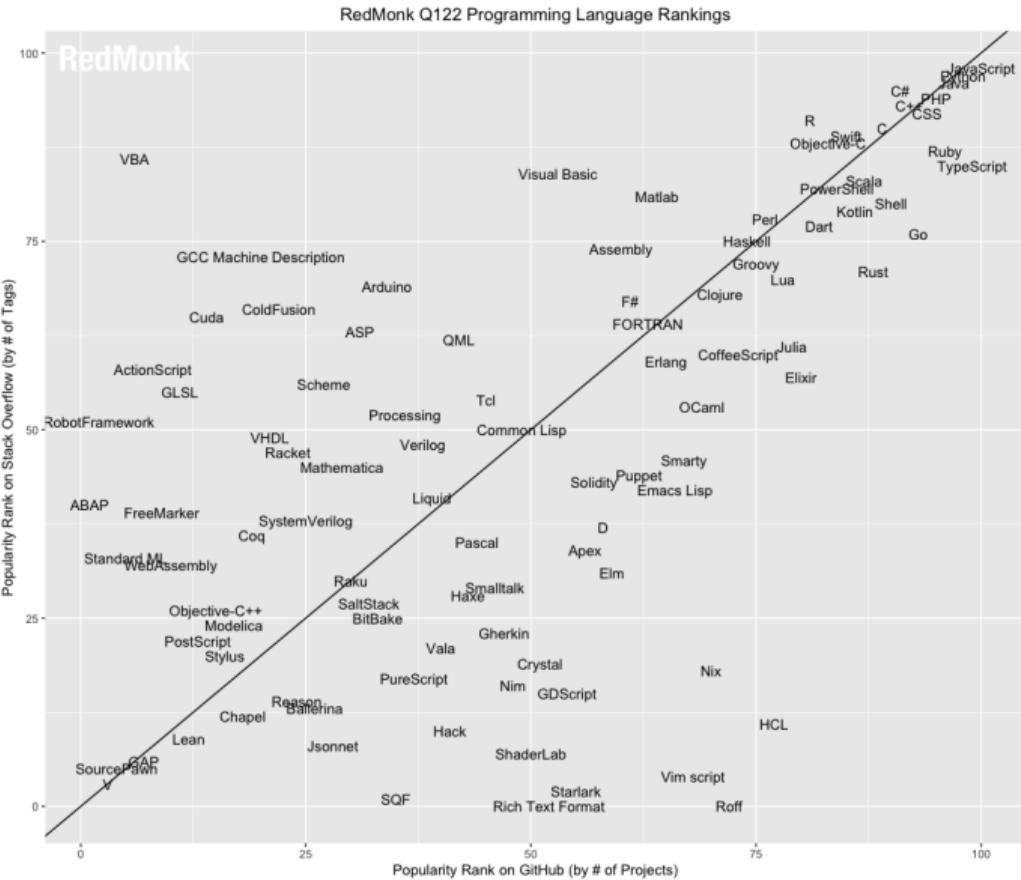


Compilation in this course:

- Start from an imperative core language
- Compile statically to a binary executable
- Classical intermediate representations and algorithms

Compilation in this course:

- Start from an imperative core language
  - Compile statically to a binary executable
  - Classical intermediate representations and algorithms
- ⇒ But language design didn't stop at C!



How to approach:

- dynamic languages (Javascript, Ruby, Python, Scheme, ...)
- hostile languages (PHP, Perl, R, Javascript, Ruby, ...)
- objects (Java, Javascript, C++, C#, ...)
- Data manipulation (SQL, GraphQL, Python for datascience, ...)
- weird features (OCaml, R, Haskell, Rust)

Enough to fill many courses!

How to approach:

- dynamic languages (Javascript, Ruby, Python, Scheme, ...)
- hostile languages (PHP, Perl, R, Javascript, Ruby, ...)
- objects (Java, Javascript, C++, C#, ...)
- Data manipulation (SQL, GraphQL, Python for datascience, ...)
- weird features (OCaml, R, Haskell, Rust)

Enough to fill many courses!

## 1 Pattern Matching Compilation

## 2 Just in Time

# Introduction: Algebraic Data Types

## Product types (a.k.a. records, tuples, structs)

```
type point = { x: int, y: int }
let distance p1 p2 =
    let dx = p2.x - p1.x in
    let dy = p2.y - p1.y in
    sqrt(dx * dx + dy * dy)
```

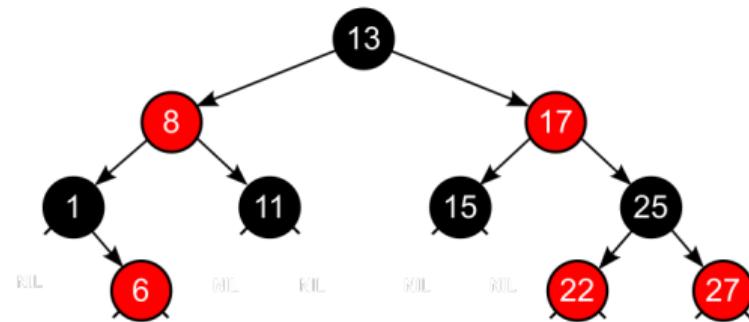
## Sum types (a.k.a. enums, tagged unions)

```
type Card = King | Queen | Jack | @Numeral of int@
let value c = match c with
    | King -> 13
    | Queen -> 12
    | Jack -> 11
    | Numeral(n) -> n
```

# Introduction: Algebraic Data Types

## Example of red-black trees

```
type color = Red | Black
type rbt =
| Empty
| Node of Color * int * RBT * RBT
```



```
let cardinal (t: RBT) : int = match t with
| Empty -> 0
| Node(_, _, t1, t2) -> 1 + cardinal(t1) + cardinal(t2)
```

# Introduction: Algebraic Data Types

## Example of red-black trees

```
type color = Red | Black
type rbt =
| Empty
| Node of Color * int * RBT * RBT
```

- Type safety, exhaustivity and non-redundancy checks
- Complex nested patterns are expressive yet concise

Rebalancing operation:

```
match c, v, t1, t2 {
| Black, z, Node(Red, y, Node(Red, x, a, b), c), d
| Black, z, Node(Red, x, a, Node(Red, y, b, c)), d
| Black, x, a, Node(Red, z, Node(Red, y, b, c), d)
| Black, x, a, Node(Red, y, b, Node(Red, z, c, d))
  -> Node(Red, y, Node(Black, x, a, b), Node(Black, z, c, d))
| a, b, c, d -> Node (a, b, c, d)
```

- How to execute this?
- How to compile this?

# Let's try in practice

Try to compile this code to if-tests:

```
let f x y z = match x,y,z with
| _, false, true -> 1
| false, true, _ -> 2
| _, _, false -> 3
| _, _, true -> 4
```

# A more complex example

$$\tau_0 = \text{None} + \text{Some}(A + B + C(u32))$$

```
match v with
| (None | Some(A)) -> 0
| Some(B) -> 1
| Some(C(n)) -> 2 + n
```

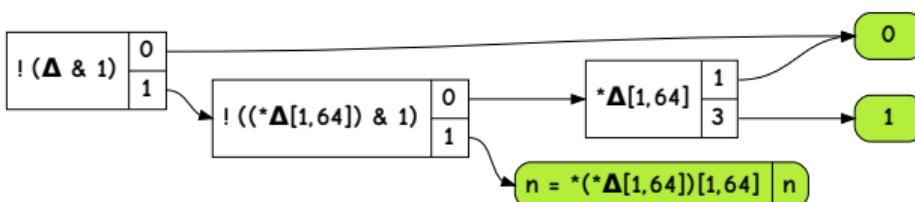
Let's try it.

# OCaml representation-specific pattern matching

## Input

```
match v with
| (None | Some(A)) -> 0
| Some(B) -> 1
| Some(C(n)) -> 2 + n
```

## Output decision tree (as graph)



## Output decision tree (as C code)

```
switch(v & 1) {
  case 0: // Some (pointer to ...)
    switch((*v)[1] & 1) {
      case 0: // C
        uint32_t n = (((*v)[1]))[1] >> 1;
        return 2 + n;
      case 1: // unit variant (A or B)
        switch((*v)[1]) {
          case 0b01: // A
            return 0;
          case 0b11: // B
            return 1;
        }
      case 1: // None: last bit is 1 (non ptr)
        return 0;
    }
}
```

# An Intermediate Representation!

- We have a pattern matching problem
  - We want a decision tree
- ⇒ We need an appropriate intermediate representation!

# The pattern Matrix

Let's try to represent a matching problem in its globality

```
match v with
| (None | Some(A)) -> 0
| Some(B) -> 1
| Some(C(n)) -> 2 + n
```

$\Rightarrow$

	.			
None   Some(A)				$\emptyset, 0$
Some(B)				$\emptyset, 1$
Some(C(n))				$\emptyset, 2$

.0	.1	.2		
-	False	True		$\emptyset, 1$
False	True	-		$\emptyset, 2$
-	-	False		$\emptyset, 3$
-	-	True		$\emptyset, 4$

```
match x,y,z with
| _, false, true -> 1
| false, true, _ -> 2
| _, _, false -> 3
| _, _, true -> 4
```

$\Rightarrow$

In the pattern matrix, columns represent “positions” in the input, and lines are patterns and outputs

# Compilation scheme

	.
<i>None</i>   <i>Some</i> ( <i>A</i> )	$\emptyset, 0$
<i>Some</i> ( <i>B</i> )	$\emptyset, 1$
<i>Some</i> ( <i>C</i> ( <i>n</i> ))	$\emptyset, 2$

- Initial pattern matrix matches the main discriminant against toplevel patterns
- Each case yields its index and an empty binding environment

# Compilation scheme

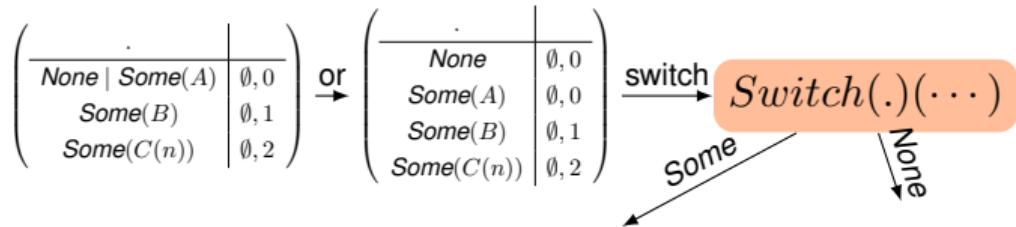
$$\left( \begin{array}{c|c} \text{None} \mid \text{Some}(A) & \emptyset, 0 \\ \text{Some}(B) & \emptyset, 1 \\ \text{Some}(C(n)) & \emptyset, 2 \end{array} \right)$$

or

$$\left( \begin{array}{c|c} \text{None} & \emptyset, 0 \\ \text{Some}(A) & \emptyset, 0 \\ \text{Some}(B) & \emptyset, 1 \\ \text{Some}(C(n)) & \emptyset, 2 \end{array} \right)$$

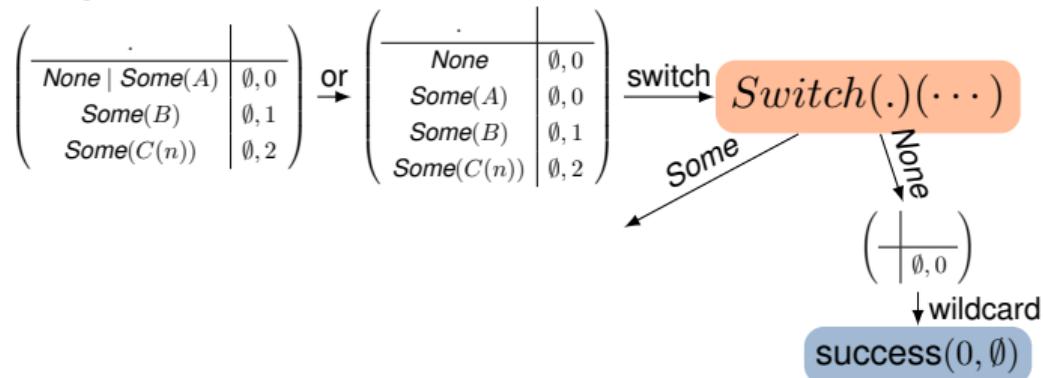
Split or-patterns

# Compilation scheme



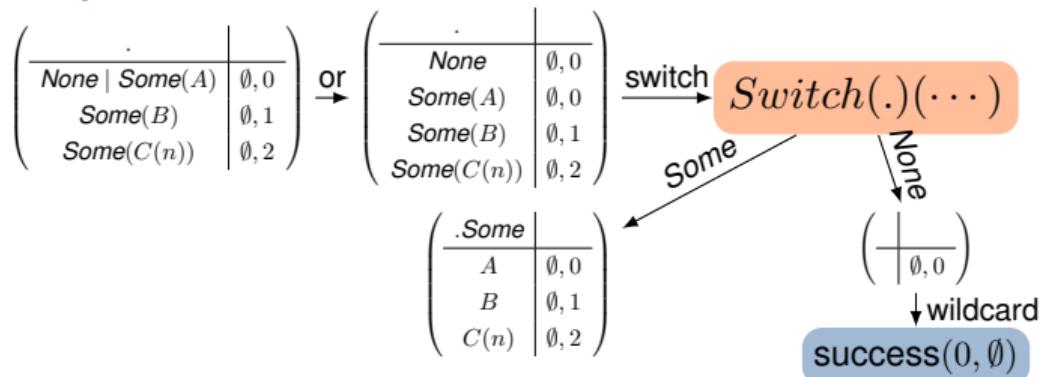
- Retrieve the head constructor of the current subterm, then branch to its associated subtree
- One branch per constructor

# Compilation scheme



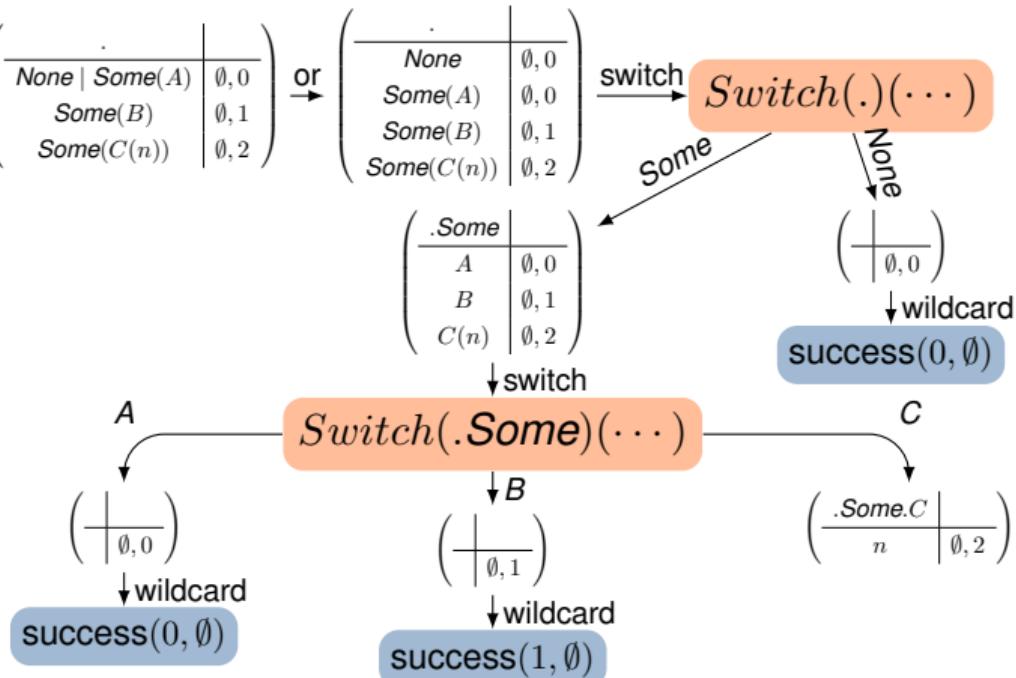
First case accepts any input

# Compilation scheme



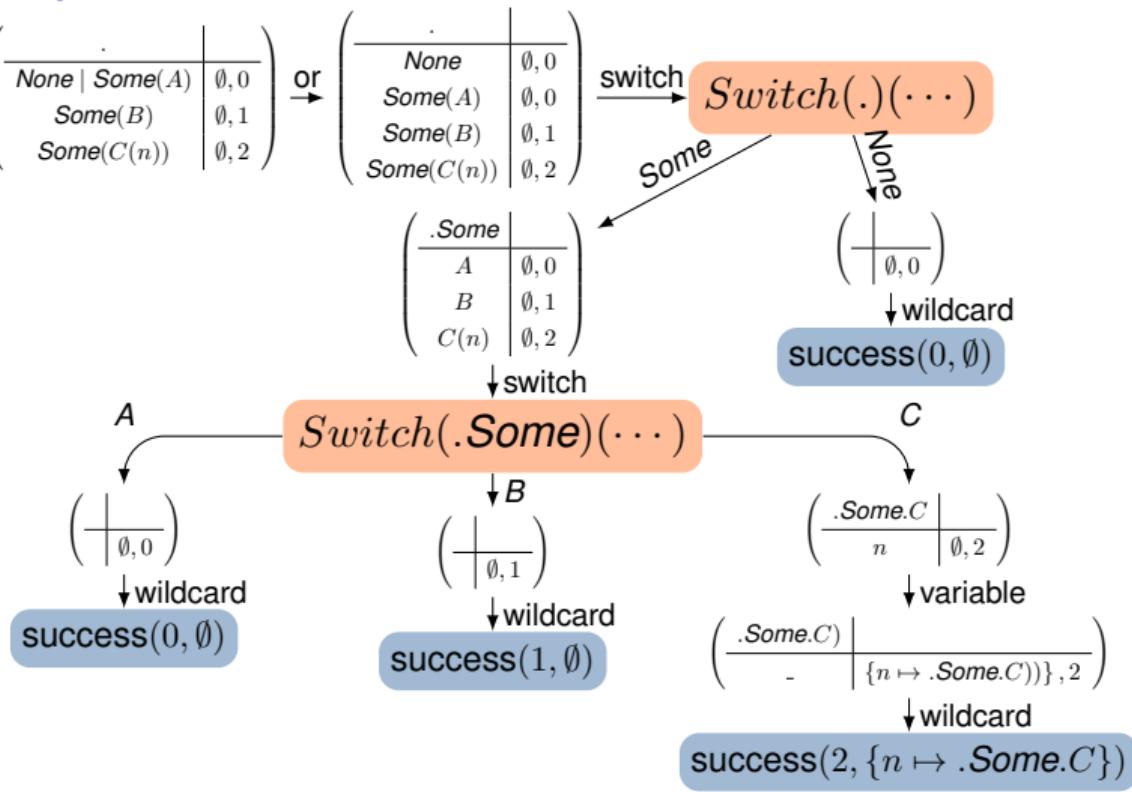
- Discard head-constructor-incompatible cases
- Focus remaining cases on the child subterm

# Compilation scheme



Inspect the current subterm

# Compilation scheme



- Bind variable pattern to the current subterm

# Compilation algorithm – Fail and Success

$\text{COMPILE}(\emptyset) = \text{unreachable}$

(No pattern case)

$$\text{COMPILE} \left( \begin{array}{ccc|c} h_1 & \dots & h_n & \\ \hline - & \dots & - & j^1, s^1 \\ \vdots & \ddots & \vdots & \vdots \\ p_1^m & \dots & p_n^m & j^m, s^m \end{array} \right) = \text{success}(j^1, s^1) \quad (\text{Wildcard case})$$

# Compilation algorithm – Variables

$$\text{COMPILE} \left( \begin{array}{ccc|c} h_1 & h_i & h_n & \\ \hline p_1^1 & p_i^1 & p_n^1 & j^1, s^1 \\ \vdots & \vdots & \vdots & \vdots \\ p_1^\ell & \cdots \boxed{x} \cdots & p_n^\ell & j^\ell, s^\ell \\ \vdots & \vdots & \vdots & \vdots \\ p_1^m & p_i^m & p_n^m & j^m, s^m \end{array} \right) = \text{COMPILE} \left( \begin{array}{ccc|c} h_1 & h_i & h_n & \\ \hline p_1^1 & \dots & p_n^1 & j^1, s^1 \\ \vdots & \vdots & \vdots & \vdots \\ p_1^\ell & \cdots \boxed{-} \cdots & p_n^\ell & j^\ell, \boxed{s^\ell \cup s_x} \\ \vdots & \vdots & \vdots & \vdots \\ p_1^m & \dots & p_n^m & j^m, s^m \end{array} \right)$$

where  $s_x = \{x \rightarrow h_i\}$

(Variable case)

# Compilation algorithm – Or

**COMPILE**

$$\left( \begin{array}{ccc|c}
 h_1 & h_i & h_n & \\
 \hline
 p_1^1 & \dots & p_n^1 & j^1, s^1 \\
 \vdots & \vdots & \vdots & \vdots \\
 p_1^\ell & \dots & p_n^\ell & j^\ell, s^\ell \\
 \vdots & \vdots & \vdots & \vdots \\
 p_1^m & \dots & p_n^m & j^m, s^m
 \end{array} \right) = \text{COMPILE} \left( \begin{array}{ccc|c}
 h_1 & h_i & h_n & \\
 \hline
 p_1^1 & \dots & p_n^1 & j^1, s^1 \\
 \vdots & \vdots & \vdots & \vdots \\
 p_1^\ell & \dots & p_n^\ell & j^\ell, s^\ell \\
 p_1^\ell & \dots & p_n^\ell & j^\ell, s^\ell \\
 \vdots & \vdots & \vdots & \vdots \\
 p_1^m & \dots & p_n^m & j^m, s^m
 \end{array} \right)$$

(Or case)

# Compilation algorithm – Switch

$\text{COMPILE}(\mathcal{P}) = \left\{ \begin{array}{l} i \leftarrow \text{PICKCOLUMN}(\mathcal{P}) \\ \mathcal{P} = \left( \begin{array}{c|ccc} h_1 & h_i & h_n & \\ \hline p_1^1 & p_i^1 & p_n^1 & j^1, s^1 \\ \vdots & \vdots & \vdots & \vdots \\ p_1^\ell & \cdots p_i^\ell \cdots & p_n^\ell & j^\ell, s^\ell \\ \vdots & \vdots & \vdots & \vdots \\ p_1^m & p_i^m & p_n^m & j^m, s^m \end{array} \right) \\ \text{Tags} = \text{GETTAGS}(p_i^1, \dots, p_i^m) \\ \forall tag \in \text{Tags}, \mathcal{P}_{tag} = \text{EXPAND}(\mathcal{P}, \text{type}(h_i), h_i) \\ \text{Switch}(h_i) \{ tag \mapsto \text{COMPILE}(\mathcal{P}_{tag}) \} \end{array} \right\}$  (Switch case)

# Compilation algorithm – Expand

Inputs		Outputs	
$\tau$	$tag$	New Headers	Matrix transformation
$\langle \tau_0, \dots, \tau_l \rangle$		$(h.0 \quad \dots \quad h.\ell)$	$\langle p_0, \dots, p_\ell \rangle \mapsto \begin{pmatrix} p_0 & \dots & p_\ell \\ - & \mapsto & - \\ - & \dots & - \end{pmatrix}$
$\sum_{1 \leq i \leq \ell} K_i(\tau_i)$	$K_{i_0}$	$(h.K_{i_0})$	$K_{i_0}(p) \mapsto \begin{pmatrix} p \end{pmatrix}$ $K_i(\dots) \mapsto \emptyset$ $- \mapsto \begin{pmatrix} - & \dots & - \end{pmatrix}$

# Compilation algorithm – Picking a column

How to pick a column?

No clear answer, we want to minimize:

(1) The longest path length, (2) The size of the decision trees.

⇒ Heuristics!

Example of heuristics:

- First row that has a pattern
- Small arity
- The most “needed” columns
- ...

# To decision trees ? in OCaml

How do we check the head constructors in reality?

# To decision trees ? in OCaml

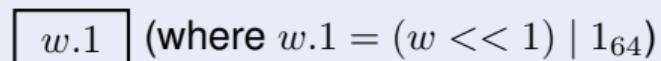
How do we check the head constructors in reality?  $\Rightarrow$  Depends on the language!

## Memory Values in OCaml

### Blocks



### Unboxed constants



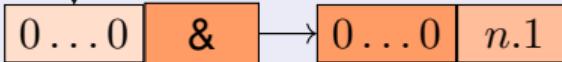
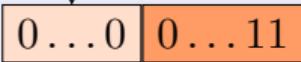
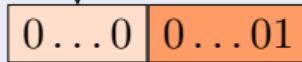
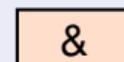
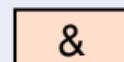
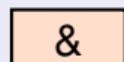
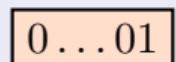
$\tau_0 = \text{None} + \text{Some}(A + B + C(u32))$  memory values

*None*

*Some ( A )*

*Some ( B )*

*Some ( C ( n ) )*



# To real decision trees!

Combine

- The base decision tree
- Specification of the language

OCaml decision tree

$\text{switch}(v \& 1)$

0  $\mapsto$   $\text{switch}((\ast v).1 \& 1)$

1  $\mapsto$   $\text{switch}((\ast v).1)$

01  $\mapsto$  0,  $\emptyset$

11  $\mapsto$  1,  $\emptyset$

0  $\mapsto$  2, { $n \mapsto (\ast(\ast v).1).1$ }

1  $\mapsto$  0,  $\emptyset$

# How to implement integer-level switches ?

The CPU doesn't have switches!

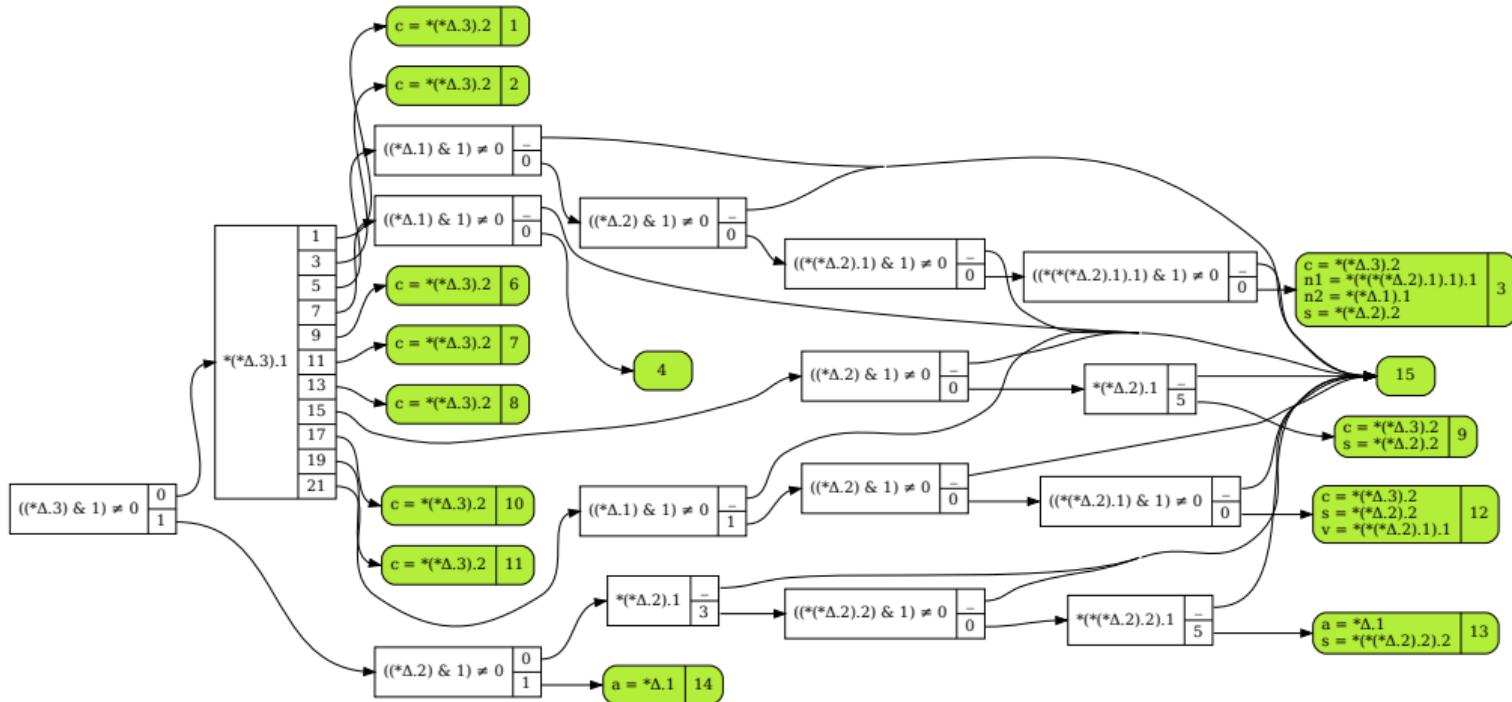
Switch implementation: highly depends on the instruction set:

- If-trees
- Bitmasks
- Jump-tables

# Big example

A big example (interpreter for a stack language):

```
matching (a, slist, clist) with
| _, _, Cons(Ldi, c) -> 1
| _, _, Cons(Push, c) -> 2
| Int(n2), Cons(Val(Int(n1)), s), Cons(I0p, c) -> 3
| Int(_), _, Cons(Test,_) -> 4
| Int(_), _, Cons(Test,_) -> 5
| _, _, Cons(Extend,c) -> 6
| _, _, Cons(Search,c) -> 7
| _, _, Cons(Pushenv,c) -> 8
| _, Cons(Env,s), Cons(Popenv,c) -> 9
| _, _, Cons(Mkclos,c) -> 10
| _, _, Cons(Mkclosrec,c) -> 11
| Clo, Cons(Val(v),s), Cons(Apply,c) -> 12
| a, Cons(Code,Cons(Env,s)), Nil -> 13
| a, Nil, Nil -> 14
| _ -> 15
```



## Exo time

Let's compile the following pattern matrix with the heuristics “first row that has a pattern”:

.	
$\langle A, B \rangle$	$\emptyset, 1$
$\langle B, - \rangle$	$\emptyset, 2$
$\langle C(x), C(y) \rangle$	$\emptyset, 3$
$\langle -, x \rangle$	$\emptyset, 4$

# Conclusion

We have seen how to compile pattern matching

- Not so trivial! Lot's of optimization opportunities
- Essential in functional languages
- Also useful elsewhere: LLVM has similar algorithms for cases on strings

Takeaway ⇒ “Niche” features deserve their compilation too

## 1 Pattern Matching Compilation

## 2 Just in Time

- Speculation
- Tracing

## A Toy Example

```
#include <stdio.h>
#include <stdlib.h>
#include <sys/mman.h>

int main(void) {
    char* program;
    int (*fnptr)(void);
    int a;
    program = mmap(NULL, 1000, PROT_EXEC | PROT_READ |
        PROT_WRITE, MAP_PRIVATE | MAP_ANONYMOUS, 0, 0);
    program[0] = 0xB8;
    program[1] = 0x34;
    program[2] = 0x12;
    program[3] = 0;
    program[4] = 0;
    program[5] = 0xC3;
    fnptr = (int (*)(void)) program;
    a = fnptr();
    printf("Result = %X\n", a);
}
```

- 1) What is the program on the left doing?
- 2) What is **this API** all about?
- 3) What does this program have to do with a just-in-time compiler?

## A Toy Example

```
#include <stdio.h>
#include <stdlib.h>
#include <sys/mman.h>

int main(void) {
    char* program;
    int (*fnptr)(void);
    int a;
    program = mmap(NULL, 1000, PROT_EXEC | PROT_READ |
    PROT_WRITE, MAP_PRIVATE | MAP_ANONYMOUS, 0, 0);
    program[0] = 0xB8;
    program[1] = 0x34;
    program[2] = 0x12;
    program[3] = 0;
    program[4] = 0;
    program[5] = 0xC3;
    fnptr = (int (*)(void)) program;
    a = fnptr();
    printf("Result = %X\n", a);
}
```

## A Toy Example

```
#include <stdio.h>
#include <stdlib.h>
#include <sys/mman.h>

int main(void) {
    char* program;
    int (*fnptr)(void);
    int a;
    program = mmap(NULL, 1000, PROT_EXEC | PROT_READ |
        PROT_WRITE, MAP_PRIVATE | MAP_ANONYMOUS, 0, 0);
    program[0] = 0xB8;
    program[1] = 0x34;
    program[2] = 0x12;
    program[3] = 0;
    program[4] = 0;
    program[5] = 0xC3;
    fnptr = (int (*)(void)) program;
    a = fnptr();
    printf("Result = %X\n", a);
}
```

## A Toy Example

```
#include <stdio.h>
#include <stdlib.h>
#include <sys/mman.h>

int main(void) {
    char* program;
    int (*fnptr)(void);
    int a;
    program = mmap(NULL, 1000, PROT_EXEC | PROT_READ |
        PROT_WRITE, MAP_PRIVATE | MAP_ANONYMOUS, 0, 0);
    program[0] = 0xB8;
    program[1] = 0x34;
    program[2] = 0x12;
    program[3] = 0;
    program[4] = 0;
    program[5] = 0xC3;
    fnptr = (int (*)(void)) program;
    a = fnptr();
    printf("Result = %X\n", a);
}
```



## Just-in-Time Compilers

- A JIT compiler translates a program into binary code while this program is being executed.
- We can compile a function as soon as it is necessary.
  - This is Google's V8 approach.
- Or we can first interpret the function, and after we realize that this function is hot, we compile it into binary.
  - This is the approach of Mozilla's IonMonkey.

1) When/where/why are just-in-time compilers usually used?

2) Can a JIT compiled program run faster than a statically compiled program?

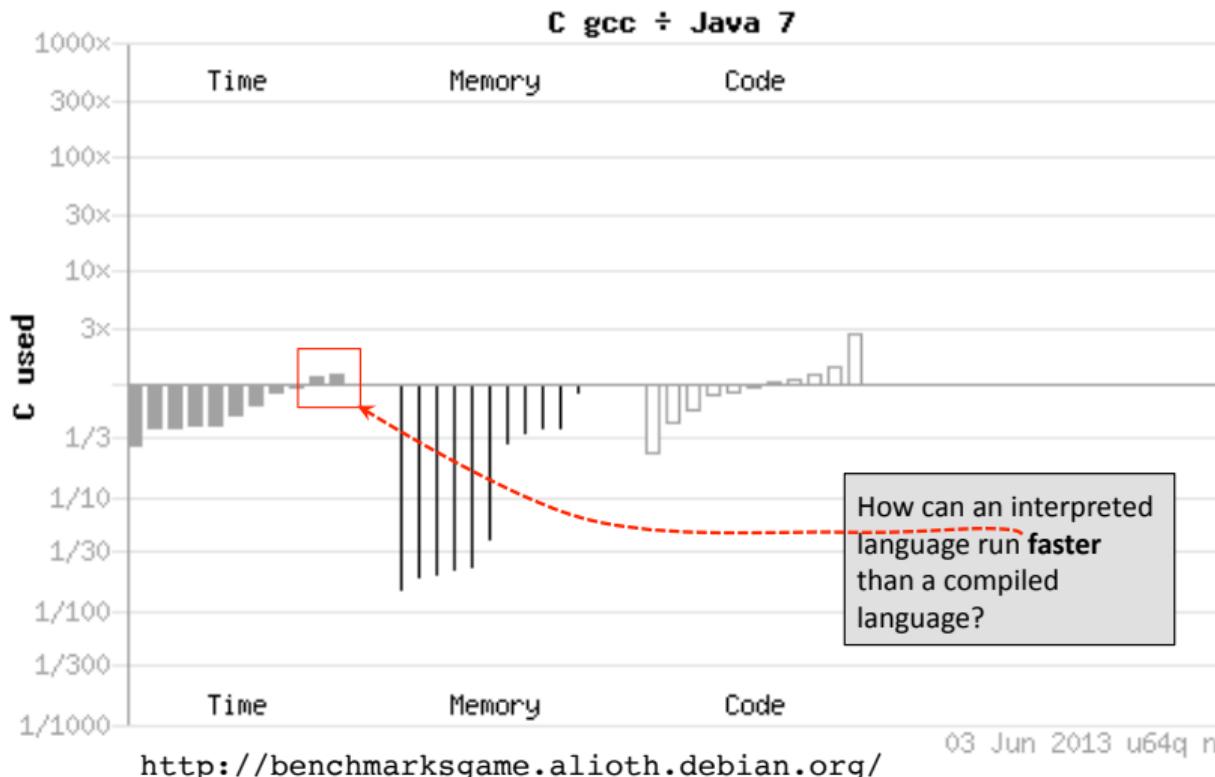
3) Which famous JIT compilers do we know?

## There are many JIT compilers around

- Java Hotspot is one of the most efficient JIT compilers in use today. It was released in 1999, and has been in use since then.
- V8 is the JavaScript JIT compiler used by Google Chrome.
- IonMonkey is the JavaScript JIT compiler used by the Mozilla Firefox.
- LuaJIT (<http://luajit.org/>) is a trace based just-in-time compiler that generates code for the Lua programming language.
- The .Net framework JITs CIL code.
- For Python we have PyPy, which runs on Cpython.



## Can JITs compete with static compilers?



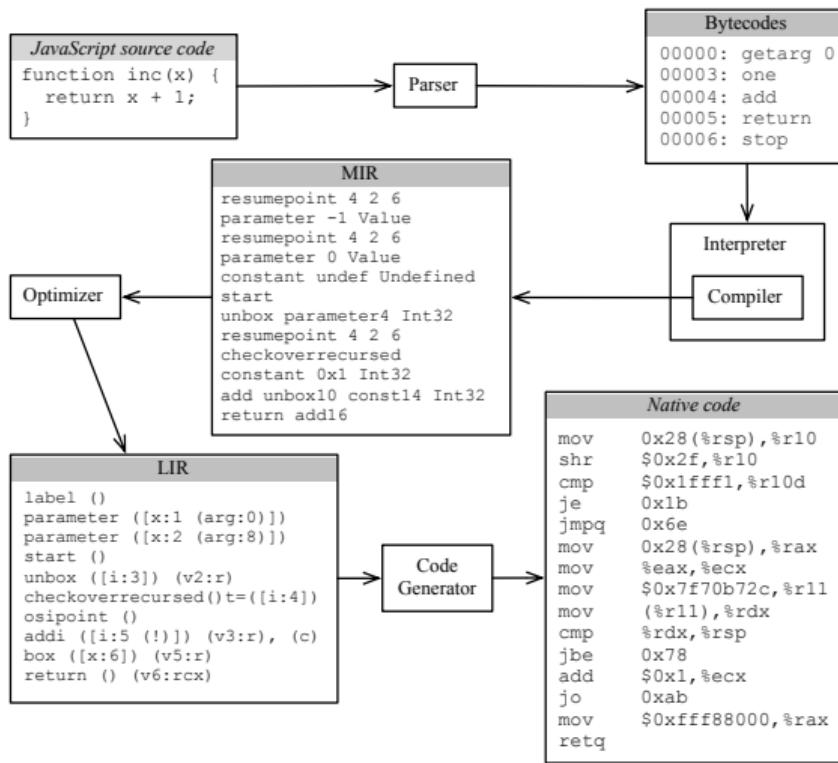


## Tradeoffs

- There are many tradeoffs involved in the JIT compilation of a program.
- The time to compile is part of the total execution time of the program.
- We may be willing to run simpler and faster optimizations (or no optimization at all) to diminish the compilation overhead.
- And we may try to look at runtime information to produce better codes.
  - Profiling is a big player here.
- The same code may be compiled many times!

Why would we compile the same code many times?

# Example: Mozilla's IonMonkey



IonMonkey is one of the JIT compilers used by the Firefox browser to execute JavaScript programs.

This compiler is tightly integrated with SpiderMonkey, the JavaScript interpreter.

SpiderMonkey invokes IonMonkey to JIT compile a function either if it is often called, or if it has a loop that executes for a long time.

Why do we have so many different intermediate representations here?

## When to Invoke the JIT Compiler?

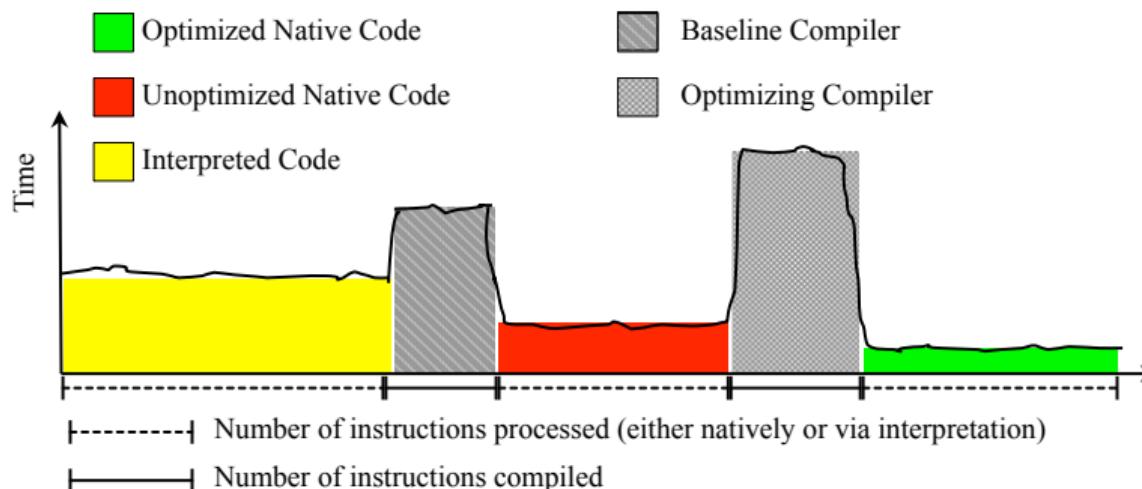
- Compilation has a cost.
  - Functions that execute only once, for a few iterations, should be interpreted.
- Compiled code runs faster.
  - Functions that are called often, or that loop for a long time should be compiled.
- And we may have different optimization levels...

How to decide when to compile a piece of code?

As an example, SpiderMonkey uses three execution modes: the first is interpretation; then we have the baseline compiler, which does not optimize the code. Finally IonMonkey kicks in, and produces highly optimized code.

## The Compilation Threshold

Many execution environments associate counters with branches. Once a counter reaches a given threshold, that code is compiled. But defining which threshold to use is very difficult, e.g., how to minimize the area of the curve below? JITs are crowded with magic numbers.



## The Million-Dollars Question

- When to invoke the JIT compiler?

- 1) Can you come up with a strategy to invoke the JIT compiler that optimizes for speed?
- 2) Do you have to execute the program a bit before calling the JIT?
- 3) How much information do you need to make a good guess?
- 4) What is the price you pay for making a wrong prediction?
- 5) Which programs are easy to predict?
- 6) Do the easy programs reflect the needs of the users?



2

## Just in Time

- Speculation
- Tracing



## Speculation

- A key trick used by JIT compilers is *speculation*.
- We may assume that a given property is true, and then we produce code that capitalizes on that speculation.
- There are many different kinds of speculation, and they are always a *gamble*:

## Speculation

- A key trick used by JIT compilers is *speculation*.
- We may assume that a given property is true, and then we produce code that capitalizes on that speculation.
- There are many different kinds of speculation, and they are always a *gamble*:
  - Let's assume that the type of a variable is an integer,
    - but if we have an integer overflow...
  - Let's assume that the properties of the object are fixed,
    - but if we add or remove something from the object...
  - Let's assume that the target of a call is always the same,
    - but if we point the function reference to another closure...





## Inline Caching

- One of the earliest, and most effective, types of specialization was *inline caching*, an optimization developed for the Smalltalk programming language<sup>⊕</sup>.
- Smalltalk is a dynamically typed programming language.
- In a nutshell, objects are represented as *hash-tables*.
- This is a very flexible programming model: we can add or remove properties of objects at will.
- Languages such as Python and Ruby also implement objects in this way.
- Today, inline caching is the key idea behind JITs's high performance when running JavaScript programs.

<sup>⊕</sup>: Efficient implementation of the smalltalk-80 system, POPL (1984)



# Virtual Tables

```

class Animal {
    public void eat() {
        System.out.println(this + " is eating");
    }
    public String toString () { return "Animal"; }
}

class Mammal extends Animal {
    public void suckMilk() {
        System.out.println(this + " is sucking");
    }
    public String toString () { return "Mammal"; }
    public void eat() {
        System.out.println(this + " is eating like a mammal");
    }
}

class Dog extends Mammal {
    public void bark() {
        System.out.println(this + " is barking");
    }
    public String toString () { return "Dog"; }
    public void eat() {
        System.out.println(this + ", is eating like a dog");
    }
}
  
```

In Java, C++, C#, and other programming languages that are compiled statically, objects contain pointers to virtual tables. Any method call can be resolved after two pointer dereferences. The first dereference finds the virtual table of the object. The second finds the target method, given a known offset.

- 1) In order for this trick to work, we need to impose some restrictions on virtual tables. Which ones?
- 2) What are the virtual tables of?
 

```

Animal a = new Animal();
Animal m = new Mammal();
Mammal d = new Dog()
      
```

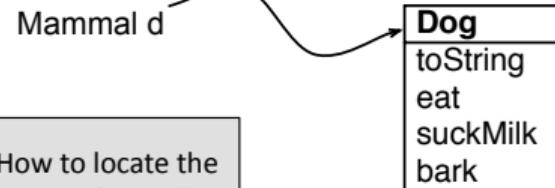
# Virtual Tables

```

class Animal {
    public void eat() {
        System.out.println(this + " is eating");
    }
    public String toString () { return "Animal"; }
}

class Mammal extends Animal {
    public void suckMilk() {
        System.out.println(this + " is sucking");
    }
    public String toString () { return "Mammal"; }
    public void eat() {
        System.out.println(this + " is eating like a mammal");
    }
}

class Dog extends Mammal {
    public void bark() {
        System.out.println(this + " is barking");
    }
    public String toString () { return "Dog"; }
    public void eat() {
        System.out.println(this + ", is eating like a dog");
    }
}
  
```



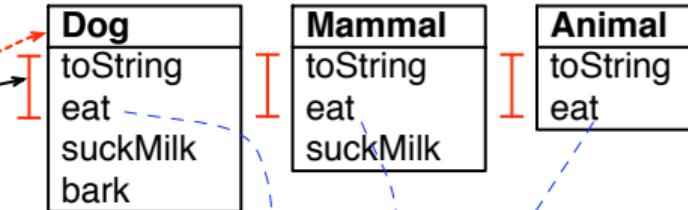
How to locate the target of d.eat()?

# Virtual Call

d.eat()

First, we need to know the table d is pointing to. This requires one pointer dereference:

Mammal d



Second, we need to know the offset of the method eat, inside the table. This offset is always the same for any class that inherits from Animal, so we can jump blindly.

Can we have virtual tables in duck typed languages?

```

public void eat() {
    System.out.println("Eats like a dog");
}

public void eat() {
    System.out.println("Eats like a mammal");
}

public void eat() {
    System.out.println("Eats like an animal");
}
  
```

# Objects in Python

```
INT_BITS = 32

def getIndex(element):
    index = element / INT_BITS
    offset = element % INT_BITS
    bit = 1 << offset
    return (index, bit)

class Set:
    def __init__(self, capacity):
        self.capacity = capacity
        self.vector = range(1+capacity/INT_BITS)
        for i in range(len(self.vector)):
            self.vector[i] = 0
    def add(self, element):
        (index, bit) = getIndex(element)
        self.vector[index] |= bit

class ErrorSet(Set):
    def checkIndex(self, element):
        if (element > self.capacity):
            raise IndexError(str(element) + " is out of range.")
    def add(self, element):
        self.checkIndex(element)
        Set.add(self, element)
        print element, "successfully added."
```

- 1) What is the program on the left doing?
- 2) What is the complexity to locate the target of a method call in Python?
- 3) Why can't calls in Python be implemented as efficiently as calls in Java?

# Using Python Objects

```
def fill(set, b, e, s):
    for i in range(b, e, s):
        set.add(i)

s0 = Set(15)
fill(s0, 10, 20, 3)

s1 = ErrorSet(15)
fill(s1, 10, 14, 3)

class X:
    def __init__(self):
        self.a = 0

fill(X(), 1, 10, 3)
>>> AttributeError: X instance
>>> has no attribute 'add'
```

- 1) What does the function fill do?
- 2) Why did the **third** call of fill failed?
- 3) What are the requirements that fill expects on its parameters?



## Duck Typing

```
def fill(set, b, e, s):
    for i in range(b, e, s):
        set.add(i)

class Num:
    def __init__(self, num):
        self.n = num
    def add(self, num):
        self.n += num
    def __str__(self):
        return str(self.n)

n = Num(3)
print n
fill(n, 1, 10, 1)
...
```

Do we get an error  
**here?**

## Duck Typing

```
class Num:  
    def __init__(self, num):  
        self.n = num  
    def add(self, num):  
        self.n += num  
    def __str__(self):  
        return str(self.n)  
  
n = Num(3)  
print n  
fill(n, 1, 10, 1)  
print n  
  
>>> 3  
>>> 48
```

The program works just fine. The only requirement that `fill` expects on its first argument is that it has a method `add` that takes two parameters. Any object that has this method, and can receive an integer on the second argument, will work with `fill`. This is called **duck typing**: *if it quacks like a duck, swims like a duck, eats like a duck, then it is a duck!*



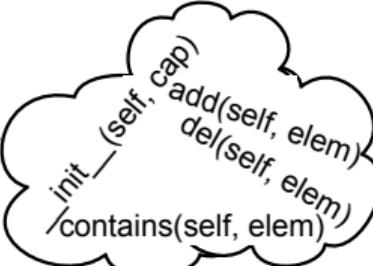
## The Price of Flexibility

- Objects, in these dynamically typed languages, are for the most part implemented as hash tables.
  - That is cool: we can add or remove methods without much hard work.
  - And mind how much code we can reuse?
- But method calls are pretty expensive.

```
def fill(set, b, e, s):  
    for i in range(b, e, s):  
        set.add(i)
```

How can we make these calls cheaper?

Mammal d



A cloud-shaped callout box containing several method names, each with a curved arrow pointing towards it from the text "Mammal d". The methods listed are: \_\_init\_\_(self, cap), add(self, elem), del(self, elem), and /contains(self, elem).

# Monomorphic Inline Cache

```
class Num:  
    def __init__(self, n):  
        self.n = n  
    def add(self, num):  
        self.n += num  
  
def fill(set, b, e, s):  
    for i in range(b, e, s):  
        set.add(i)  
  
n = Num(3)  
print n  
fill(n, 1, 10, 1)  
print n  
  
=>> 3  
=>> 48
```

The first time we generate code for a call, we can check the target of that call. We know this target, because we are generating code at runtime!

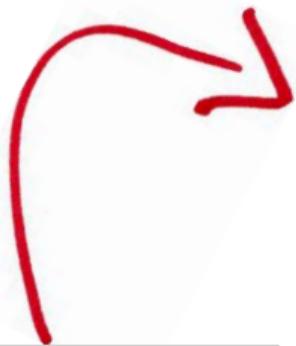
```
fill(set, b, e, s):  
    for i in range(b, e, s):  
        if isinstance(set, Num):  
            set.n += i  
        else:  
            add = lookup(set, "add")  
            add(set, i)
```

Could you optimize this code even further using classic compiler transformations?

## Inlining on the Method

- We can also speculate on the method name, instead of doing it on the calling site:

```
fill(set, b, e, s):
    for i in range(b, e, s):
        if isinstance(set, Num):
            set.n += i
        else:
            add = lookup(set, "add")
            add(set, i)
```



```
fill(set, b, e, s):
    for i in range(b, e, s):
        __f_add(set, i)

__f_add(o, e):
    if isinstance(o, Num):
        o.n += e
    else:
        f = lookup(o, "add")
        f(o, e)
```

- 1) Is there any advantage to this approach, when compared to inlining at the call site?
- 2) Is there any disadvantage?
- 3) Which one is likely to change more often?



## Polymorphic Calls

- If the target of a call changes during the execution of the program, then we have a polymorphic call.
- A monomorphic inline cache would have to be invalidated, and we would fall back into the expensive quest.

```
>>> l = [Num(1), Set(1), Num(2), Set(2), Num(3), Set(3)]  
>>> for o in l:  
...     o.add(2)  
...
```

Is there anything we could do to optimize this case?

# Polymorphic Calls

- If the target of a call changes during the execution of the program, then we have a polymorphic call.
- A monomorphic inline cache would have to be invalidated, and we would fall back into the expensive quest.

```
>>> l = [Num(1), Set(1), Num(2), Set(2),
       Num(3), Set(3)]
>>> for o in l:
...     o.add(2)
...
```

Would it not be better just to have the code below?

```
for i in range(b, e, s):
    lookup(set, "add")
    ...

```

```
fill(set, b, e, s):
    for i in range(b, e, s):
        __f_add(set, i)

__f_add(o, e):
    if isinstance(o, Num):
        o.n = e
    elif isinstance(o, Set):
        (index, bit) = getIndex(e)
        o.vector[index] |= bit
    else:
        f = lookup(o, "add")
        f(o, e)
```

# The Speculative Nature of Inline Caching

- Python – as well as JavaScript, Ruby, Lua and other very dynamic languages – allows the user to add or remove methods from an object.
- If such changes in the layout of an object happen, then the representation of that object must be recompiled. In this case, we need to update the inline cache.

```
from Set import INT_BITS, getIndex, Set

def errorAdd(self, element):
    if (element > self.capacity):
        raise IndexError(str(element) +
                         " is out of range.")
    else:
        (index, bit) = getIndex(element)
        self.vector[index] |= bit
        print element, "added successfully!"

Set.add = errorAdd
s = Set(60)
s.errorAdd(59)
s.remove(59)
```

## The Benefits of the Inline Cache

These numbers have been obtained by Ahn *et al.* for JavaScript, in the Chrome V8 compiler<sup>◊</sup>:

- Monomorphic inline cache hit:
  - 10 instructions
- Polymorphic Inline cache hit:
  - 35 instructions if there are 10 types
  - 60 instructions if there are 20 types
- Inline cache miss: 1,000 – 4,000 instructions.

Which factors could justify these numbers?

<sup>◊</sup>: Improving JavaScript Performance by Deconstructing the Type System, PLDI (2014)

2

## Just in Time

- Speculation
- Tracing



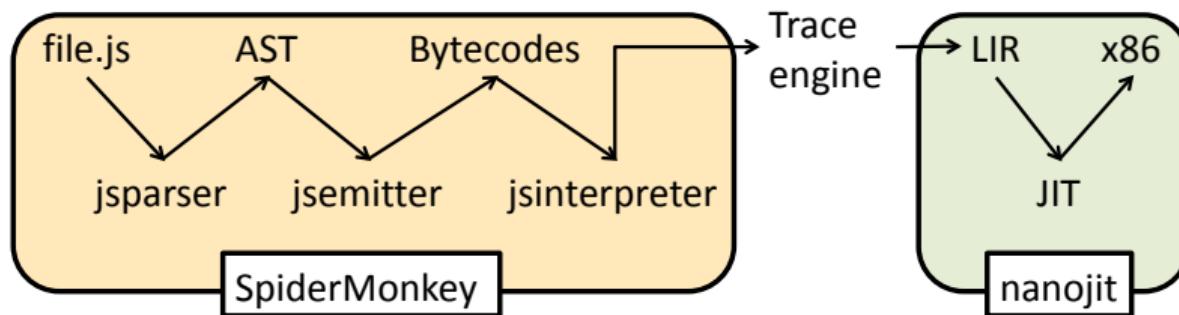
## What is a JIT trace compiler?

- A trace-based JIT compiler translates only the most executed paths in the program's control flow to machine code.
- A trace is a linear sequence of code, that represents a hot path in the program.
- Two or more traces can be combined into a tree.
- Execution alternates between traces and interpreter.

What are the advantages and disadvantages of trace compilation over traditional method compilation?

## The anatomy of a trace compiler

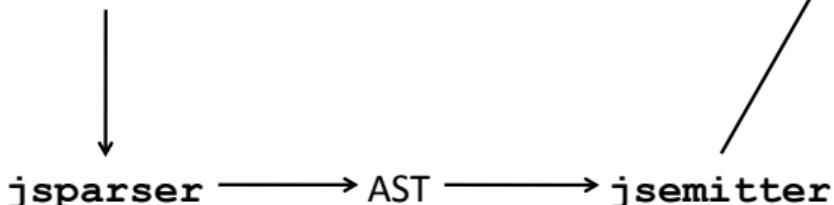
- TraceMonkey is the trace based JIT compiler used in the Mozilla Firefox Browser.



## From source to bytecodes

```
function foo(n) {
  var sum = 0;
  for(i = 0; i < n; i++) {
    sum+=i;
  }
  return sum;
}
```

Can you see the correspondence between source code and bytecodes?



00: **getname n**  
 02: **setlocal 0**  
 06: **zero**  
 07: **setlocal 1**  
 11: **zero**  
 12: **setlocal 2**  
 16: **goto 35 (19)**  
 19: **trace**  
 20: **getlocal 1**  
 23: **getlocal 2**  
 26: **add**  
 27: **setlocal 1**  
 31: **localinc 2**  
 35: **getlocal 2**  
 38: **getlocal 0**  
 41: **lt**  
 42: **ifne 19 (-23)**  
 45: **getlocal 1**  
 48: **return**  
 49: **stop**

## The trace engine kicks in

- TraceMonkey interprets the bytecodes.
- Once a **loop** is found, it may decide to ask Nanojit to transform it into machine code (e.g., x86, ARM).
  - Nanojit reads LIR and produces x86
- Hence, TraceMonkey must convert **this** trace of bytecodes into LIR



```
00: getname n
02: setlocal 0
06: zero
07: setlocal 1
11: zero
12: setlocal 2
16: goto 35 (19)
19: trace
20: getlocal 1
23: getlocal 2
26: add
27: setlocal 1
31: localinc 2
35: getlocal 2
38: getlocal 0
41: lt
42: ifne 19 (-23)
45: getlocal 1
48: return
49: stop
```

## Bytecodes

```

00: getname n
02: setlocal 0
06: zero
07: setlocal 1
11: zero
12: setlocal 2
16: goto 35 (19)
19: trace
20: getlocal 1
23: getlocal 2
26: add
27: setlocal 1
31: localinc 2
35: getlocal 2
38: getlocal 0
41: lt
42: ifne 19 (-23)
45: getlocal 1
48: return
49: stop
  
```

## From bytecodes to LIR

```

L: load "i" %r1
    load "sum" %r2
    add %r1 %r2 %r1
    %p0 = ovf()
    bra %p0 Exit1
    store %r1 "sum"
    inc %r2
    store %r2 "i"
    %p0 = ovf()
    bra %p0 Exit2
    load "i" %r0
    load "n" %r1
    lt %p0 %r0 %r1
    bne %p0 L
  
```

Nanojit LIR

## Bytecodes

```

00: getname n
02: setlocal 0
06: zero
07: setlocal 1
11: zero
12: setlocal 2
16: goto 35 (19)

```

```

19: trace
20: getlocal 1
23: getlocal 2
26: add
27: setlocal 1
31: localinc 2
35: getlocal 2
38: getlocal 0
41: lt
42: ifne 19 (-23)

```

```

45: getlocal 1
48: return
49: stop

```

## From bytecodes to LIR

```

L: load "i" %r1
load "sum" %r2
add %r1 %r2 %r1


%p0 = ovf()


bra %p0 Exit1
store %r1 "sum"
inc %r2
store %r2 "i"


%p0 = ovf()


bra %p0 Exit2
load "i" %r0
load "n" %r1
lt %p0 %r0 %r1
bne %p0 L

```

**Overflow tests** are required by operations such as add, sub, inc, dec, mul.

Nanojit LIR

## Why do we have overflow tests?

- Many scripting languages represent numbers as floating-point values.
  - Arithmetic operations are not very efficient.
- The compiler sometimes is able to infer that these numbers can be used as integers.
  - But floating-point numbers are larger than integers.
  - This is another example of speculative optimization.
  - Thus, every arithmetic operation that might cause an overflow must be preceded by a test. If the test fails, then the runtime engine must change the number's type back to floating-point.

## From LIR to assembly

```

L: load "i" %r1
   load "sum" %r2
   add %r1 %r2 %r1
   %p0 = ovf()
   bra %p0 Exit1
   store %r1 "sum"
   inc %r2
   store %r2 "i"
   %p0 = ovf()
   bra %p0 Exit2
   load "i" %r0
   load "n" %r1
   lt %p0 %r0 %r1
   bne %p0 L
  
```

Nanojit LIR

 **The overflow tests** are also translated into machine code.

Can you come up with an optimization to eliminate some of the overflow checks?

## x86 Assembly

```

L: movl -32(%ebp), %eax
   movl %eax, -20(%ebp)
   movl -28(%ebp), %eax
   movl %eax, -16(%ebp)
   movl -20(%ebp), %edx
   leal -16(%ebp), %eax
   addl %edx, (%eax)

   call _ovf
   testl %eax, %eax
   jne Exit1

   movl -16(%ebp), %eax
   movl %eax, -28(%ebp)
   leal -20(%ebp), %eax
   incl (%eax)

   call _ovf
   testl %eax, %eax
   jne Exit2

   movl -24(%ebp), %eax
   movl %eax, -12(%ebp)
   movl -20(%ebp), %eax
   cmpl -12(%ebp), %eax
   jl L
  
```

## How to eliminate the redundant tests

- We use range analysis:
  - Find the range of integer values that a variable might hold during the execution of the trace.

```
function foo(n) {  
    var sum = 0;  
    for(i = 0; i < n; i++) {  
        sum+=i;  
    }  
    return sum;  
}
```

**Example:** if we know that n is 10, and i is always less than n, then we will never have an overflow if we add 1 to i.

## Cheating at runtime

- A **static analysis** must be very conservative: if we do not know for sure the value of n, then we must assume that it may be anywhere in  $[-\infty, +\infty]$ .
- **However, we are not a static analysis!**

```
function foo(n) {  
    var sum = 0;  
    for(i = 0; i < n; i++) {  
        sum+=i;  
    }  
    return sum;  
}
```

- We are compiling at runtime!
- To know the value of n, just ask the interpreter.



## How the algorithm works

- Create a constraint graph.
  - While the trace is translated to LIR.
- Propagate range intervals.
  - Before sending the LIR to Nanojit.
  - Using infinite precision arithmetic.
- Eliminate tests whenever it is safe to do so.
  - We tell Nanojit that code for some overflow tests should not be produced.

## The constraint graph

- We have four categories of vertices:

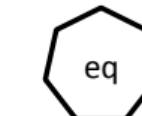
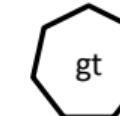
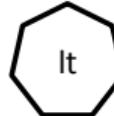
Name



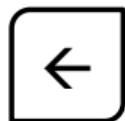
Arithmetic



Relational



Assignment



## Building the constraint graph

**19: trace**

**20: getlocal sum**

**23: getlocal i**

**26: add**

**27: setlocal sum**

**31: localinc i**

**35: getlocal n**

**38: getlocal i**

**41: lt**

**42: ifne 19 (-23)**

- We start building the constraint graph once TraceMonkey starts recording a trace.
- TraceMonkey starts at the branch instruction, which is the first instruction visited in the loop.
  - Although it is at the end of the trace.

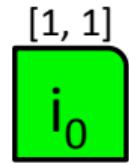
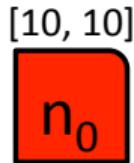
In terms of code generation,  
can you recognize the pattern  
of bytecodes created for the  
test if  $n < i$  goto L?

19: trace  
20: getlocal *sum*  
23: getlocal *i*  
26: add  
27: setlocal *sum*  
31: localinc  
35: getlocal *n*   
38: getlocal *i*  
41: It  
42: ifne 19 (-23)

[10, 10]  
  
*n<sub>0</sub>*



We check the interpreter's stack to find out that *n* is 10.

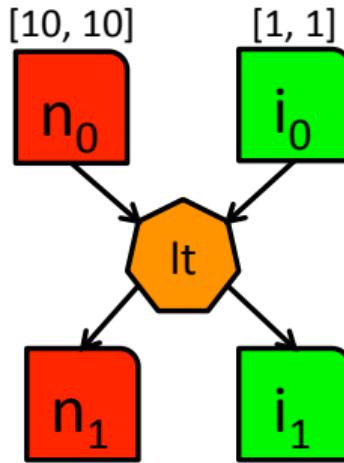


- 19: trace
- 20: getlocal *sum*
- 23: getlocal *i*
- 26: add
- 27: setlocal *sum*
- 31: localinc *i*
- 35: getlocal *i* 
- 38: getlocal *i*
- 41: It
- 42: ifne 19 (-23)

*i* started holding 0, but at this point, its value is already 1.

Why we do not initialize *i* with 0 in our constraint graph?

- 19: trace
- 20: getlocal sum
- 23: getlocal i
- 26: add
- 27: setlocal sum
- 31: localinc i
- 35: getlocal n
- 38: getlocal i
- 41: It
- 42: ifne 19 (-23)

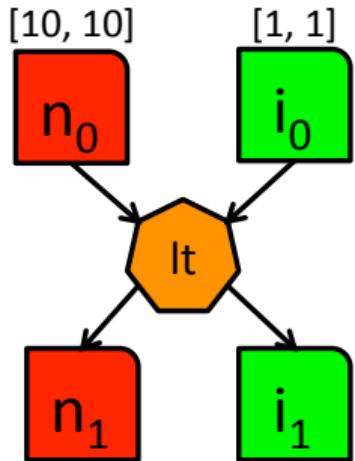


Comparisons are great! We can learn new information about variables

Now we know that i is less than 10, so we rename it to  $i_1$   
 we also rename n

We are renaming after conditionals.  
 Which program representation are we creating *dynamically*?

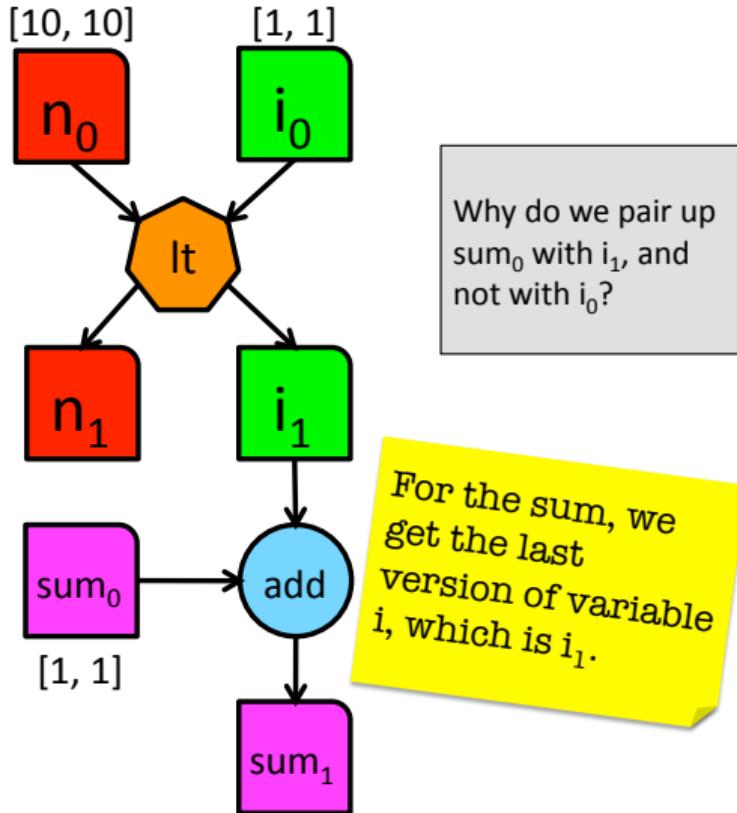
19: trace  
20: getlocal sum  
23: getlocal i  
26: add  
27: setlocal sum  
31: localinc i  
35: getlocal n  
38: getlocal i  
41: It  
42: ifne 19 (-23)



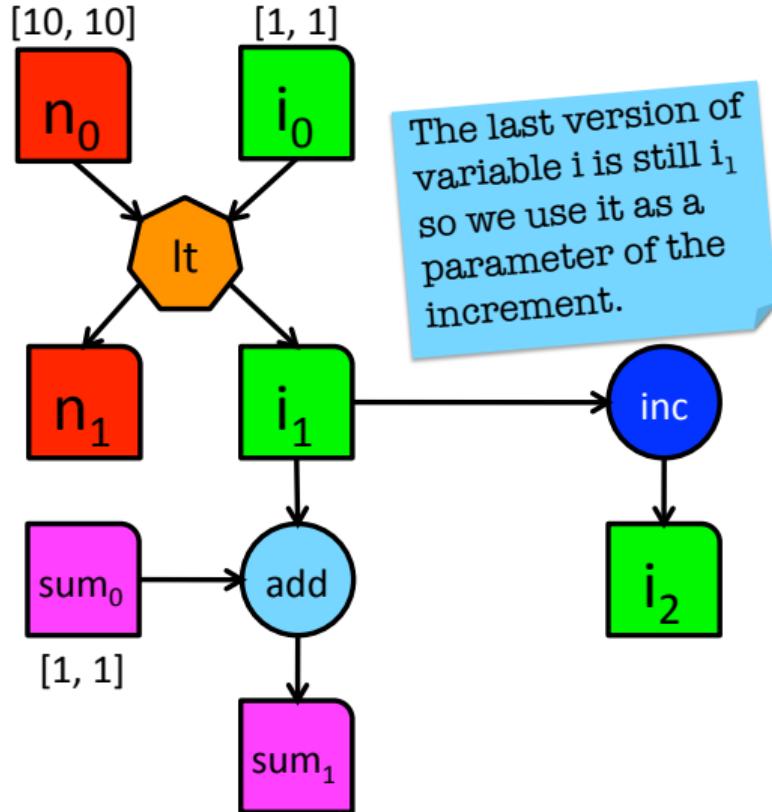
sum<sub>0</sub>  
[1, 1]

The current value of sum on the interpreter's stack is 1.

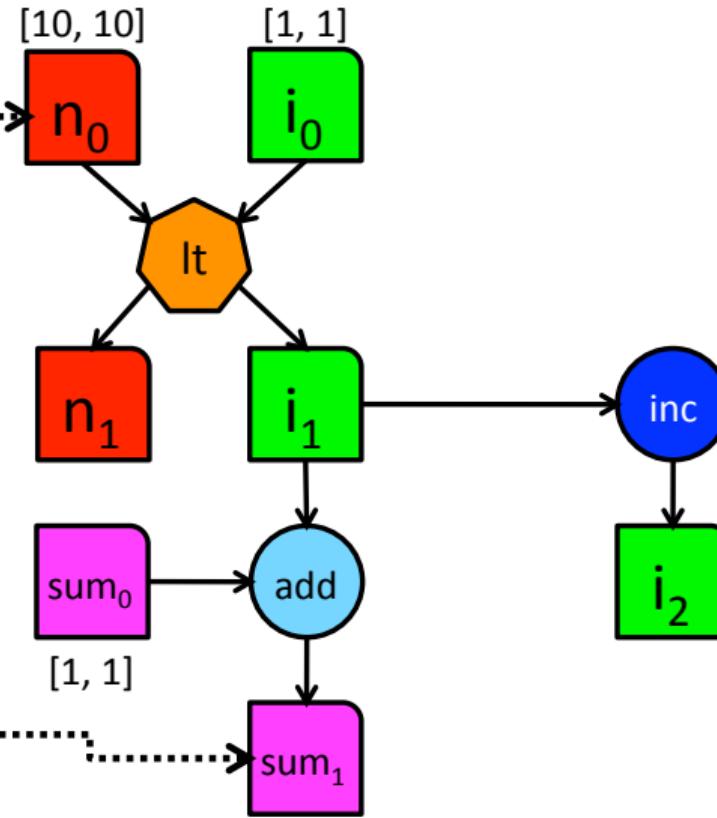
19: trace  
 20: getlocal sum  
 23: getlocal i  
 26: add ←  
 27: setlocal sum  
 31: localinc i  
 35: getlocal n  
 38: getlocal i  
 41: It  
 42: ifne 19 (-23)



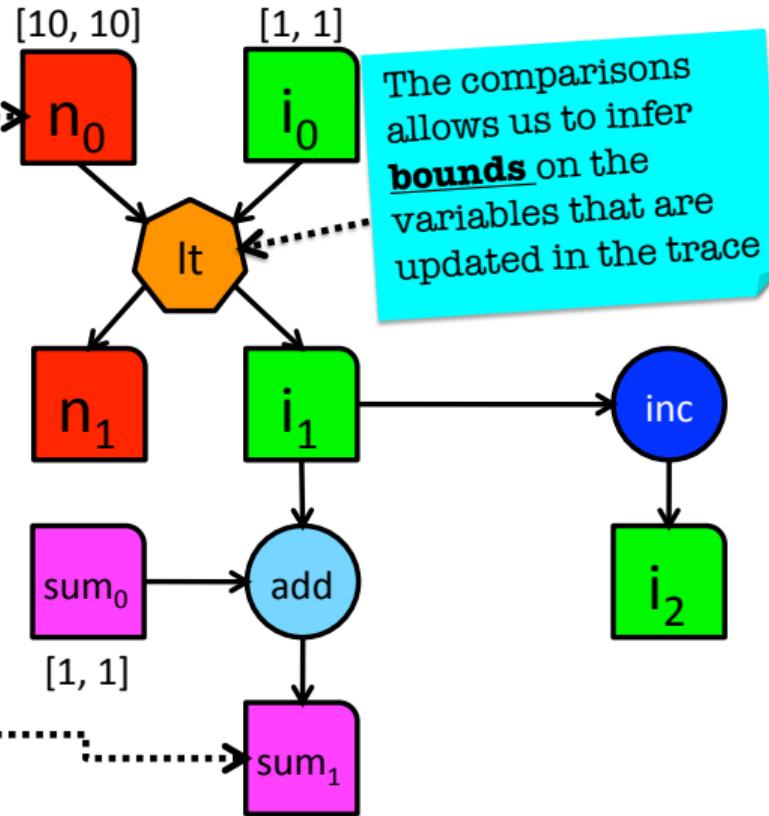
19: trace  
 20: getlocal sum  
 23: getlocal i  
 26: add  
 27: setlocal sum  
 31: localinc i  
 35: getlocal n  
 38: getlocal i  
 41: It  
 42: ifne 19 (-23)



Some variables, like  $n$ , have not been updated inside the trace. We know that they are **constants**.

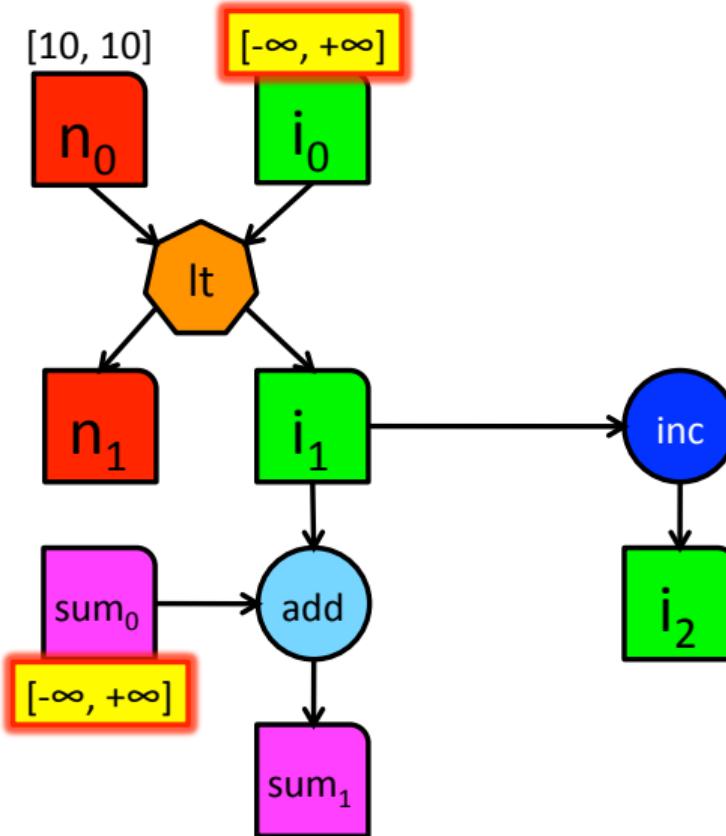


Some variables, like  $n$ , have not been updated inside the trace. We know that they are **constants**.



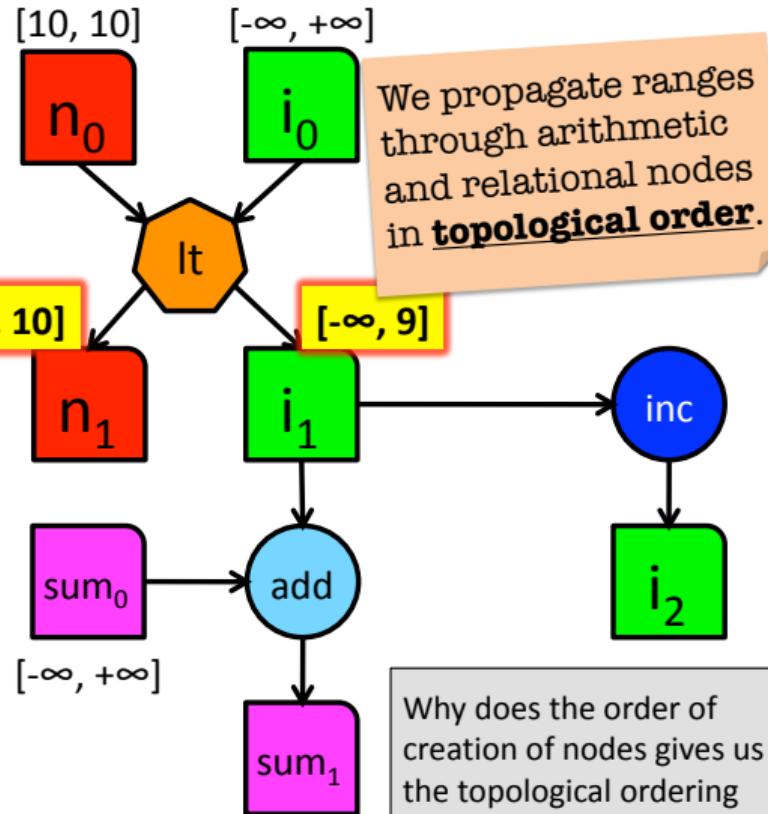
The next phase of our algorithm is the propagation of range intervals.

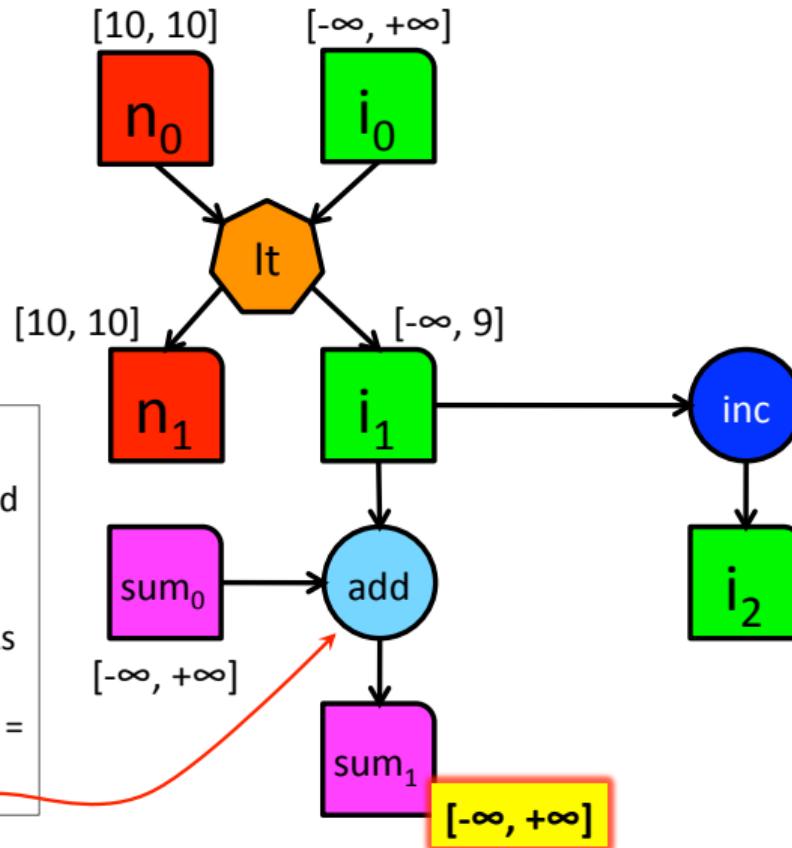
We start by assigning **conservative** i.e,  $[-\infty, +\infty]$ , bounds to the ranges of variables updated inside the trace.



19: trace  
 20: getlocal sum  
 23: getlocal i  
**26: add**  
 27: setlocal sum  
**31: localinc i**  
 35: getlocal n  
 38: getlocal i  
**41: lt**  
 42: ifne 19 (-23)

2  
3  
1

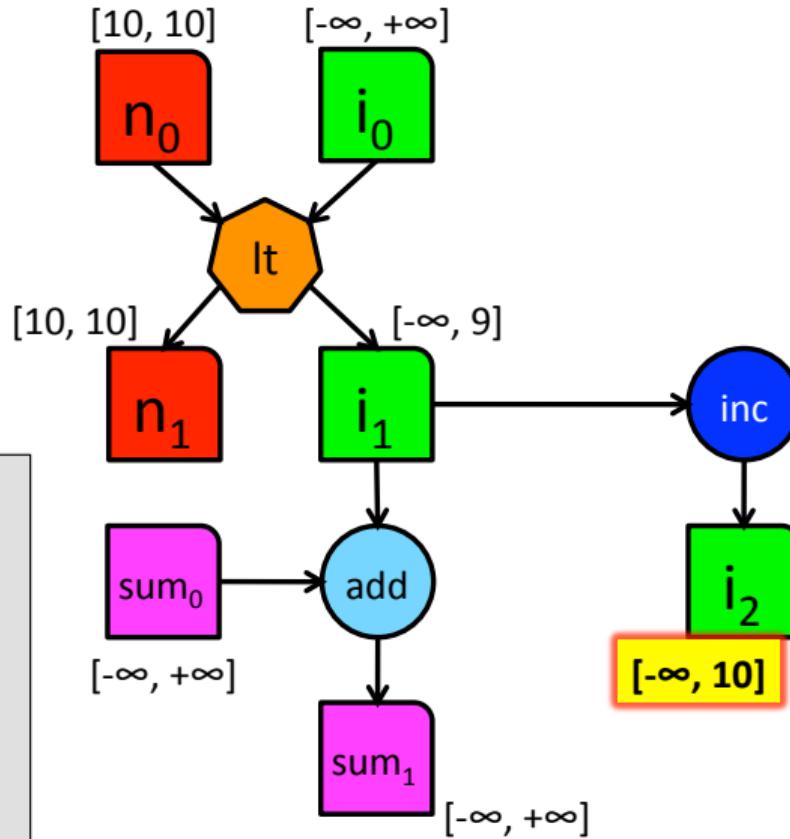




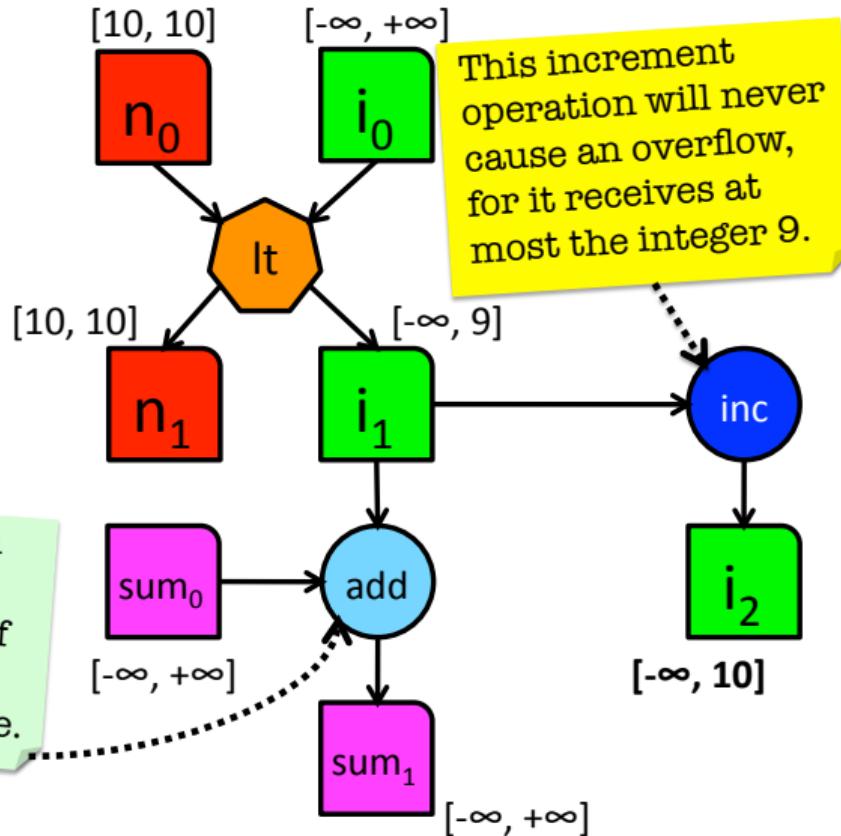
We are doing abstract interpretation. We need an abstract semantics for every operation in our constraint graph. As an example, we know that  $[-\infty, +\infty] + [-\infty, 9] = [-\infty, +\infty]$

After range propagation we can check which overflow tests are really necessary.

- 1) How many overflow checks do we have in this constraint graph?
- 2) Is there any overflow check that is redundant?

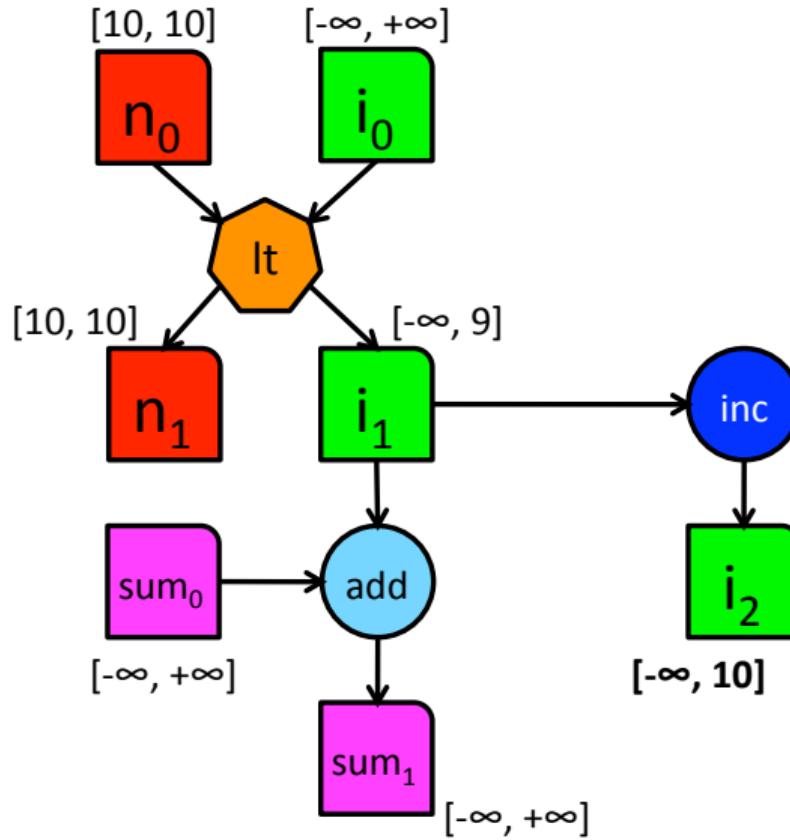


After range propagation we can check which overflow tests are really necessary.



```

L: load "i" %r1
   load "sum" %r2
   add %r1 %r2 %r1
   sp0 = ovf()
   bra %p0 Exit1
   store %r1 "sum"
   inc %r2
   store %r2 "i"
   sp0 = ovf()
   bra %p0 Exit2
   load "i" %r0
   load "n" %r1
   lt %p0 %r0 %r1
   bne %p0 L
  
```



## JIT – Conclusion

Just-In-Time compilation combines dynamic runtime information with static compilation

Practical approach: Use whatever available to make it fast.

The most used compilers in the world are Web Browsers!

## Further in Compilation

Many other “compilations” than what we have seen: dynamic languages; objects, functional, and other paradigms; Data manipulation, and other stranger computations mode (see DM from previous years!).

Recent “fun” example: implementation of the French Tax System.

# Domain Specific Languages: the new frontier

Different domains have very different computations: meteo simulations, genome analysis, cryptography, 3D rendering, Machine Learning, ...

⇒ The current frontier: How to provide nice languages and efficient compilations for these varied use-cases?