

Reusing just-in-time compiled code



Jan Vitek
Northeastern University, Boston
Czech Technical University, Prague

Thanks for having me here
I have too many slides
Work done in US, CZ, FR & IN with
students and collaborators

How can JavaScript be that fast?

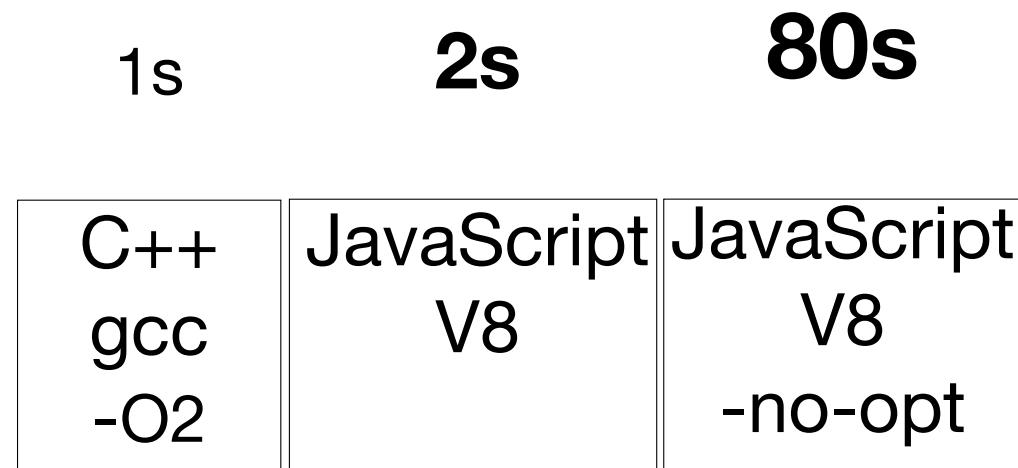


as bad as Python,
as fast as C,
how?

JIT my dear

```
function sorted(x) {  
    for (var i=1; i<x.length; i++)  
        if (x[i] < x[i-1]) return 0;  
    return 1;  
}
```

```
function sorted(x) {
    for (var i=1; i<x.length; i++)
        if (x[i] < x[i-1]) return 0;
    return 1;
}
```



measurements obtained in 2019
called with a large vector of floats

```

function sorted(x) {
  for (var i=1; i<x.length; i++)
    if (x[i] < x[i-1]) return 0;
  return 1;
}

function `[]`(x,i) {
  if (typeof(x) != array) error()
  if (typeof(i) != int)
    i = convert(i, int)
  return get(x, i)
}
function `-(`(a,b) {
  if (typeof(a) == float) {
    if (typeof(b) != float)
      b = convert(b, float)
    return subf(a.val, b.val)
  }
  if (typeof(b) == int) {
    ...
  ...
}

```

```

function `<`(a,b) {
  if (typeof(a) == float) {
    if (typeof(b) != float)
      b = convert(b, float)
    return ltf(a.val, b.val)
  }
  if (typeof(b) == int) {
    ...
  ...
}

```

Code is highly polymorphic,
but what if we could specialize to float arrays?

How can JavaScript be that fast?

```
function sorted(x) {  
    for (var i=1; i<x.length; i++)  
        t1 = get(x, i)  
        t2 = subi(i,1)  
        t3 = get(x,t2)  
        t4 = t1.val  
        t5 = t3.val  
        t6 = ltf(t4,t5)  
        if (t6) return 0  
    return 1  
}
```

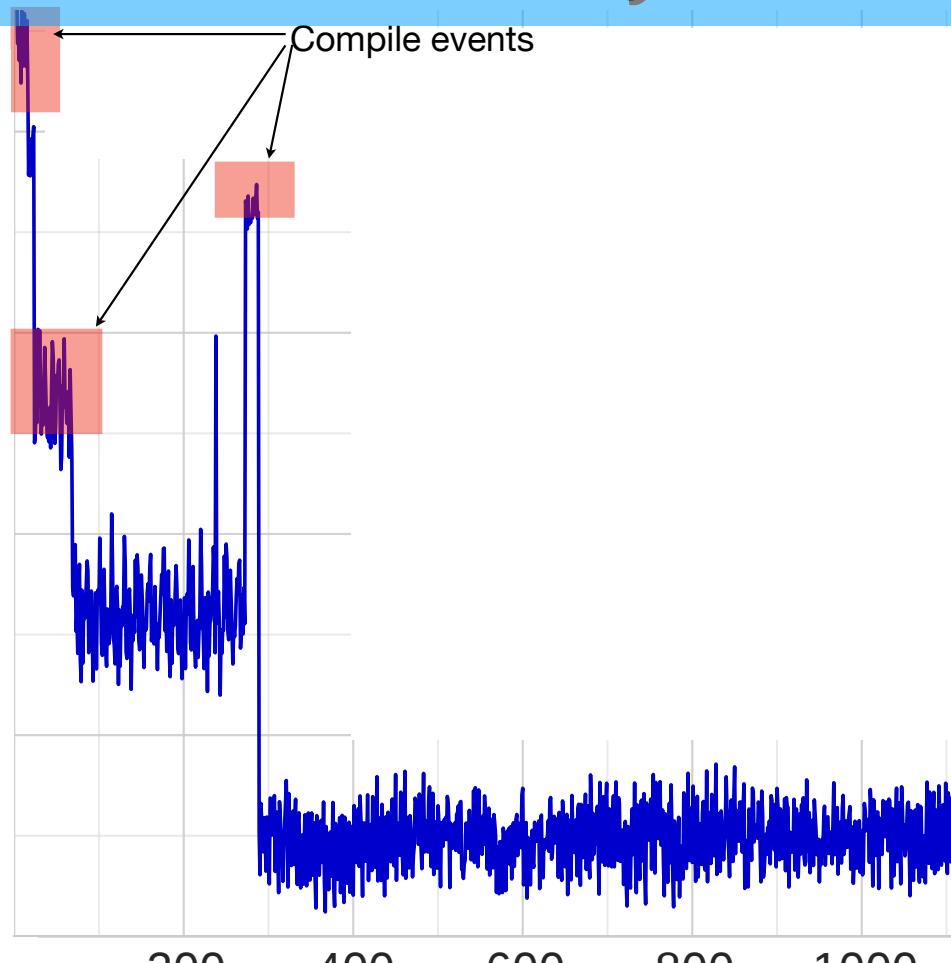
Faster version, but semantically incorrect

How can JavaScript be that fast?

```
function sorted(x) {  
  DeoptIf(typeof x != floatarray)  
  for (var i=1; i<x.length; i++)  
    t1 = get(x, i)  
    t2 = subi(i,1)  
    t3 = get(x,t2)  
    t4 = t1.val  
    t5 = t3.val  
    t6 = ltf(t4,t5)  
    if (t6) return 0  
  return 1  
}
```

correct, but requires on-stack-replacement (OSR)
if argument misspeculation occurs

How can JavaScript be that fast?



A brief history of just in time compilation

[Mitchell J.G. 1970]

Lazy compilation of uncommon branches (aka TracingJIT)

[Deutsch, Shiffman 1984]

Inline caches speculate on method receiver class

[Chambers, Ungar 1989]

Specialization based on receiver types

[Holzle, Chamber, Ungar 1992]

Dynamic deoptimization of optimized code

Speculations



Formally Verified Speculation and Deoptimization in a JIT Compiler

AURÈLE BARRIÈRE, Univ Rennes, Inria, CNRS, IRISA, France

SANDRINE BLAZY, Univ Rennes, Inria, CNRS, IRISA, France

OLIVIER FLÜCKIGER, Northeastern University, USA

DAVID PICHARDIE, Univ Rennes, Inria, CNRS, IRISA, France

JAN VITEK, Northeastern University / Czech Technical University, USA

[POPL21]

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

[POPL19]

```
function f(x, y, z)
version default:
    r1=1
L1:
    r2=x*y
    if z==7
        r3=x*x
        r4=y*y
        return r1+r3+r4
    else
        return r2
```

```
function f(x, y, z)
version optimized:
    L4:
        anchor f.default.L1=[x, y, z, r1=1]
        r2=y*y
        assume [x=1, z=75] ⚓ L4
        return 5626+r2
version default:
    ...
    ...
```

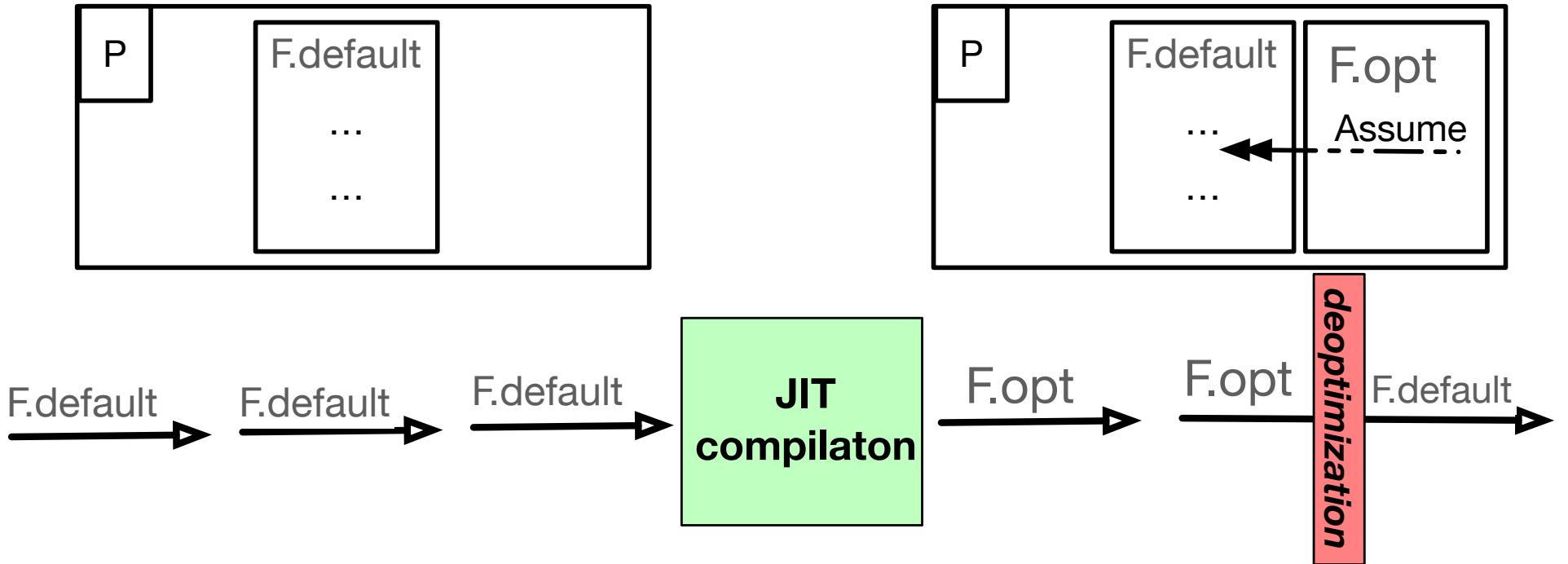
each function can be compiled to multiple ‘versions’

```
function f(x, y, z)
version default:
    r1=1
L1:
    r2=x*y
    if z==7
        r3=x*x
        r4=y*y
        return r1+r3+r4
    else
        return r2
```

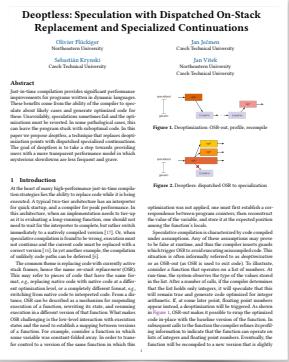
```
function f(x, y, z)
version optimized:
    L4:
        anchor f.default.L1=[x, y, z, r1=1]
        r2=y*y
        assume [x=1, z=75] ⚓ L4
        return 5626+r2
version default:
```

...

deoptimization is reified into:
anchor = how to reconstruct the stack
assume = guard condition



deoptimization jumps into less-optimized function (OSR-out),
we could also jump into more optimized function (OSR-in)



[Deoptless: Speculation with Dispatched On-Stack Replacement and Specialized Continuations, PLDI'22]

Contexts



Contextual Dispatch for Function Specialization

OLIVIER FLÜCKIGER, Northeastern University, USA

GUIDO CHARI, ASAPP INC, Argentina

MING-HO YEE, Northeastern University, USA

JAN JEČMEN, Czech Technical University, Czechia

JAKOB HAIN, Northeastern University, USA

JAN VITEK, Northeastern University, USA and Czech Technical University, Czechia

[OOPSLA 20]

how to reduce count of run-time deopts?
more versions, fewer assumes

```
x = read()
```

```
y = x[[1]] as int
```

```
z = max(x) +
```

```
max(y, 0)
```

Different uses of same function
are likely to deopt

```
function max(a, b=a, warn=FALSE)
```

```
version default:
```

```
if (warn) print(...)  
return (b < a) ? a : b
```

```
version v1 ( $\top, \perp, \perp$ ):
```

```
return a
```

```
version v2 (int[1], float[1],  $\perp$ ):
```

```
r1=a[[1]] as float
```

```
r2=b[[1]]
```

```
return ltf(r2,r1) ? a : b
```

```
x = read()
```

```
y = x[[1]] as int
```

```
z = max(x) +
```

```
max(y, 0)
```

Contextual dispatch uses
best version of a function
for each call site

```
function max(a, b=a, warn=FALSE)  
version default:
```

```
| if (warn) print(...)  
| return (b < a) ? a : b
```

```
version v1 ( $\top, \perp, \perp$ ):
```

```
| return a
```

```
version v2 (int[1], float[1],  $\perp$ ):
```

```
| r1=a[[1]] as float
```

```
| r2=b[[1]]
```

```
| return ltf(r2, r1) ? a : b
```

```
x = read()
```

```
y = x[[1]] as int
```

```
z = max(x) +
```

```
max(y, 0)
```

Contexts are constructed
at run-time based on static
and dynamic information

```
function max(a, b=a, warn=FALSE)
version default:
| if (warn) print(...)
| return (b < a) ? a : b

version v1 (T, L, L):
| return a

version v2 (int[1], float[1], L):
| r1=a[[1]] as float
| r2=b[[1]]
| return ltf(r2,r1) ? a : b
```

Reuse



Reusing JITed Code

* Jan Vitek

Reusing Just-in-Time Compiled Code

ANONYMOUS AUTHOR(S)

Most code is executed more than once. If not entire programs then, at least, libraries often remain unchanged from one program run to the next. Just-in-time compilers expend considerable effort gathering insights about code they compiled many times before, and often end up generating the same binary over and over again. This paper explores how to reuse compiled code across runs of different programs to cut down on the warm-up costs of dynamic languages. We present an approach that uses *speculative contextual dispatch* to select versions of functions from an *off-line curated code repository*. That repository is a persistent database of previously compiled functions indexed by the context under which they were compiled. Performance can be improved by curating the repository off-line to remove redundant code and optimize dispatch. We assess practicality by extending R̄, a compiler for the R programming language, and evaluating its performance over a set of benchmarks and real-world programs.

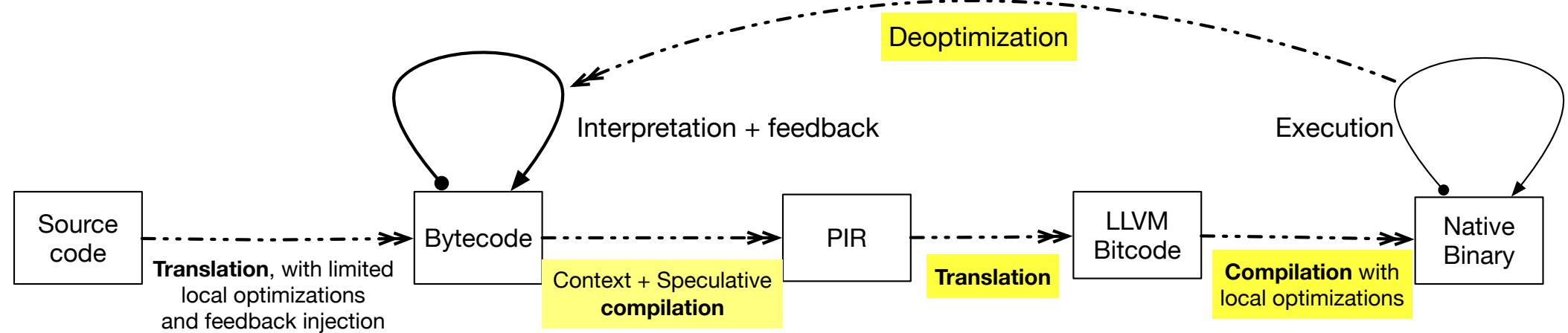
CCS Concepts: • Software and its engineering → Just-in-time compilers; Dynamic analysis;

Additional Key Words and Phrases: Specialization, Code reuse.

[XXX 23]

* IFIP * April 2023 * act 4 * Reuse

Trouble in Paradise?



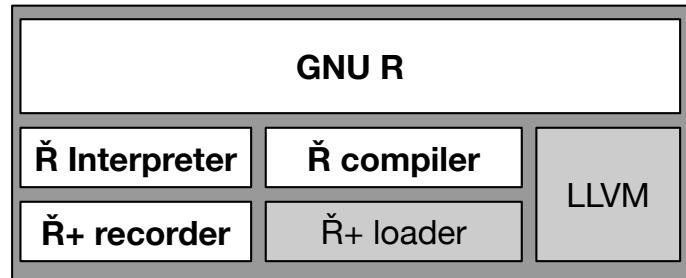
JIT warmup times can be high.
JITs have no memory across runs,
learns same truths over and over again.

[R Melts Brains: An IR for First-Class Environments and Lazy Effectful Arguments, DLS19]

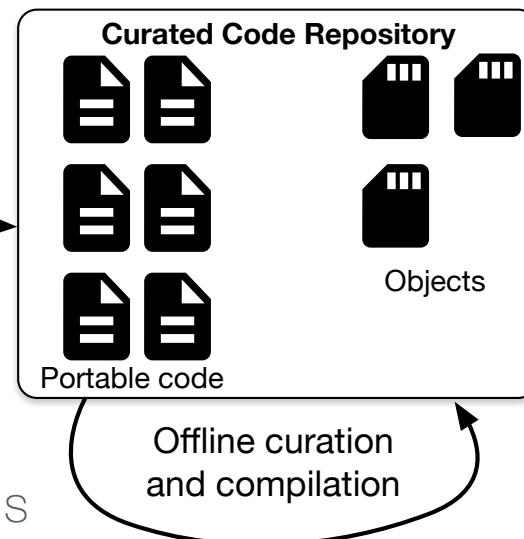


Ŕ+ allows JIT record and replay

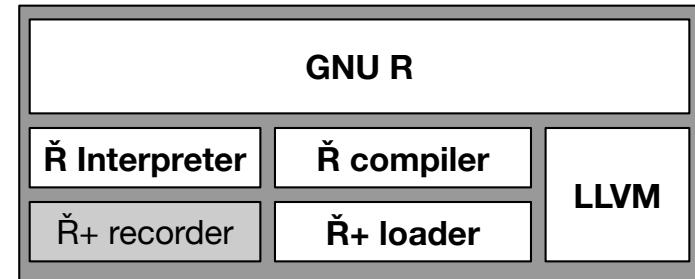
Ŕ+ in recording mode



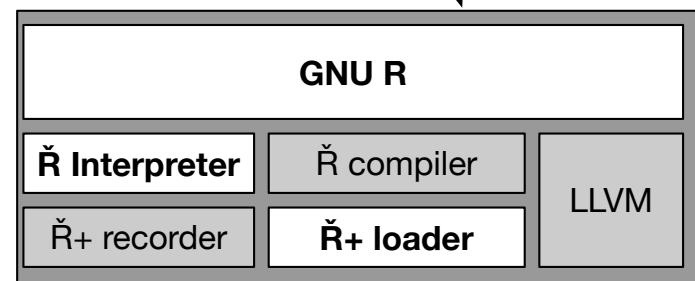
Recording output portable
intermediate code files
(LLVM bitcode and metadata)



Ŕ+ in replay mode



Replay loads executables
(native code and metadata)



Many interesting engineering details
in making code reusable but only
one interesting scientific questions

When is it safe to reuse a compiled function?

$\text{compile}(f) \rightarrow V_1 \dots \text{compile}(f) \rightarrow V_2$

Compilations of same source code may yield different binaries.

What is the soundness criteria for reuse?

```
function f(a) {  
    print(a+global)  
}
```

Let's optimize a simple function
operating over integers and strings
and reading a global variable

```
function f(a)
version default:
    r1 = global
L1:
    r2 = a
    r3 = &add
    r4 = r3(r2, r1)
    print(r4)
```

The default version with local
variables runs interpreted

```
function f(a)
version default:
r1 = global
L1:
r2 = a
r3 = &add
r4 = r3(r2, r1)
print(r4)
```

```
global = 42
f(42)
```

[r1 , r2 , r3 , r4]
[int , int , ADD , int]

Interpreter records feedback
including type approximations,
branches and function pointers

```
function f(a)
version default:
r1 = global
L1:
r2 = a
r3 = &add
r4 = r3(r2, r1)
print(r4)
```

```
global = 42
f(42)
```

[r1 , r2 , r3 , r4]
[int , int , ADD , int]

```
version V1(int):
r1 = global
L2:
anchor f.default.L1=[a, r1=r1]
r2 = a
assume [r1:int] ⚓ L2
r3 = &add
r4 = r3(r2, r1)
print(r4)
```

```
function f(a)
version default:
r1 = global
L1:
r2 = a
r3 = &add
r4 = r3(r2, r1)
print(r4)
```

```
global = 42
f(42)
```

```
[ r1 , r2 , r3 , r4 ]
[ int , int , ADD , int ]
```

```
version V2(int):
r1 = global
L2:
anchor f.default.L1=[a, r1=r1]
r2 = a
assume [r1:int] ⚓ L2
r3 = &add
assume [r3 = 0xADD] ⚓ L2
r4 = iadd(r2, r1)
print(r4)
```

```
function f(a)
version default:
r1 = global
L1:
r2 = a
r3 = &add
r4 = r3(r2, r1)
print(r4)
```

```
global = 42
f(42)
f("42")
```

```
[ r1 , r2 , r3 , r4 ]
[ int , str , ADD , int ]
```

```
version V3(int|str):
r1 = global
L2:
anchor c=[a, r1=r1]
r2 = a
assume [r1:int, r2~int] ⚓ L2
r3 = convert_int(r2)
r4 = &add
r5 = r4(a, r1)
print(r5)
```

```
(int)[int,int,ADD,int]
version V1:
    r1 = global
L2:
    anchor f.default.L1=[a, r1=r1]
    r2 = a
    assume [r1:int] ⚡ L2
    r3 = &add
    r4 = r3(r2, r1)
    print(r4)
```

```
(int)[int,int,ADD,int]
version V2:
    r1 = global
L2:
    anchor f.default.L1=[a, r1=r1]
    r2 = a
    assume [r1:int] ⚡ L2
    r3 = &add
    assume [r3 = 0xADD] ⚡ L2
    r4 = iadd(r2, r1)
    print(r4)
```

```
(int|str)[int,int|str,ADD,int]
version V3(int|str):
    r1 = global
L2:
    anchor c=[a, r1=r1]
    r2 = a
    assume [r1:int, r2~int] ⚡ L2
    r3 = convert_int(r2)
    r4 = &add
    r5 = r4(a, r1)
    print(r5)
```

$(int)[int,int,ADD,int] = (int)[int,int,ADD,int] < (int\str)[int,int\str,ADD,int]$

version V1:

assume [r1:int] ⚓ L2

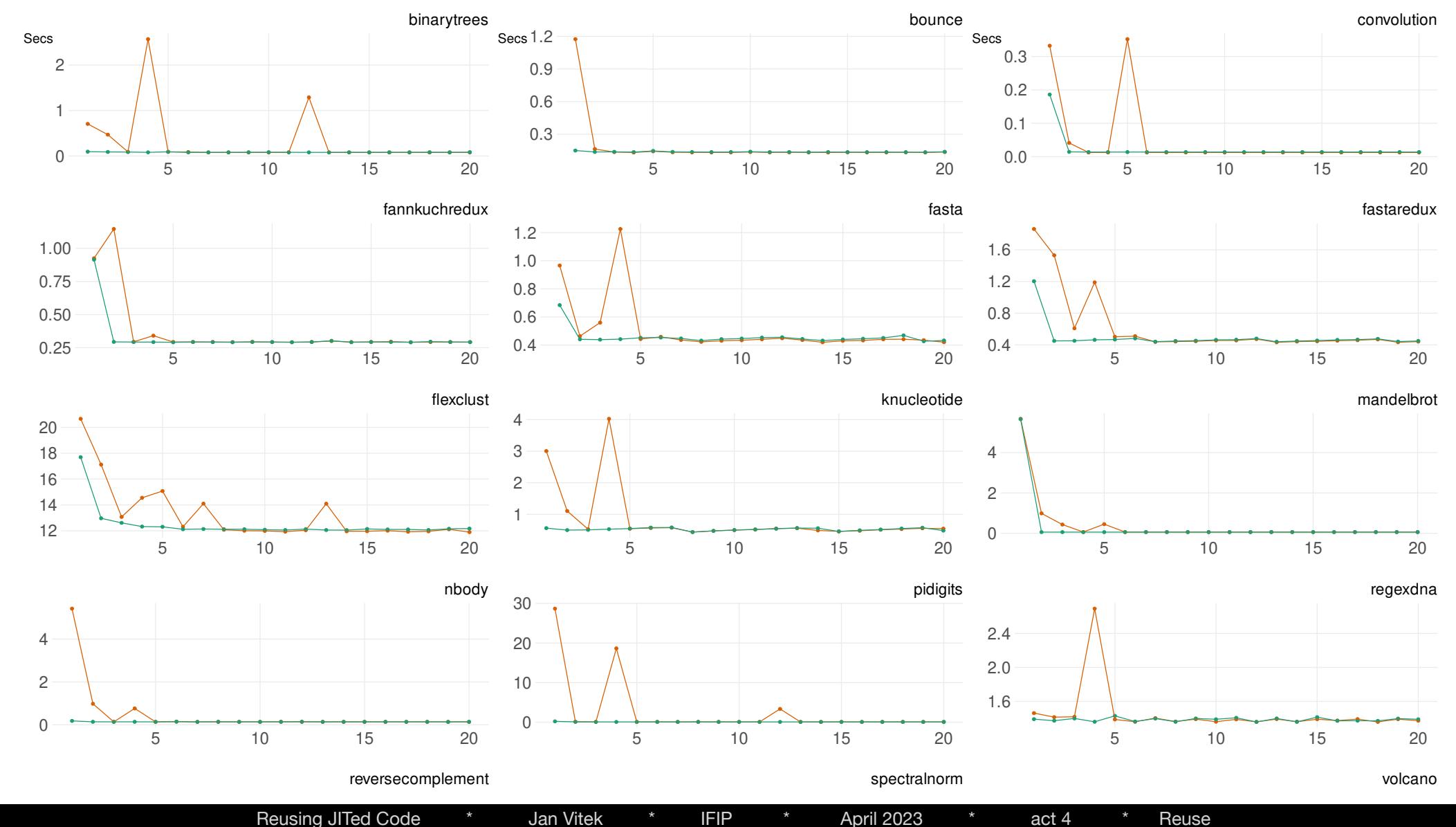
version V2:

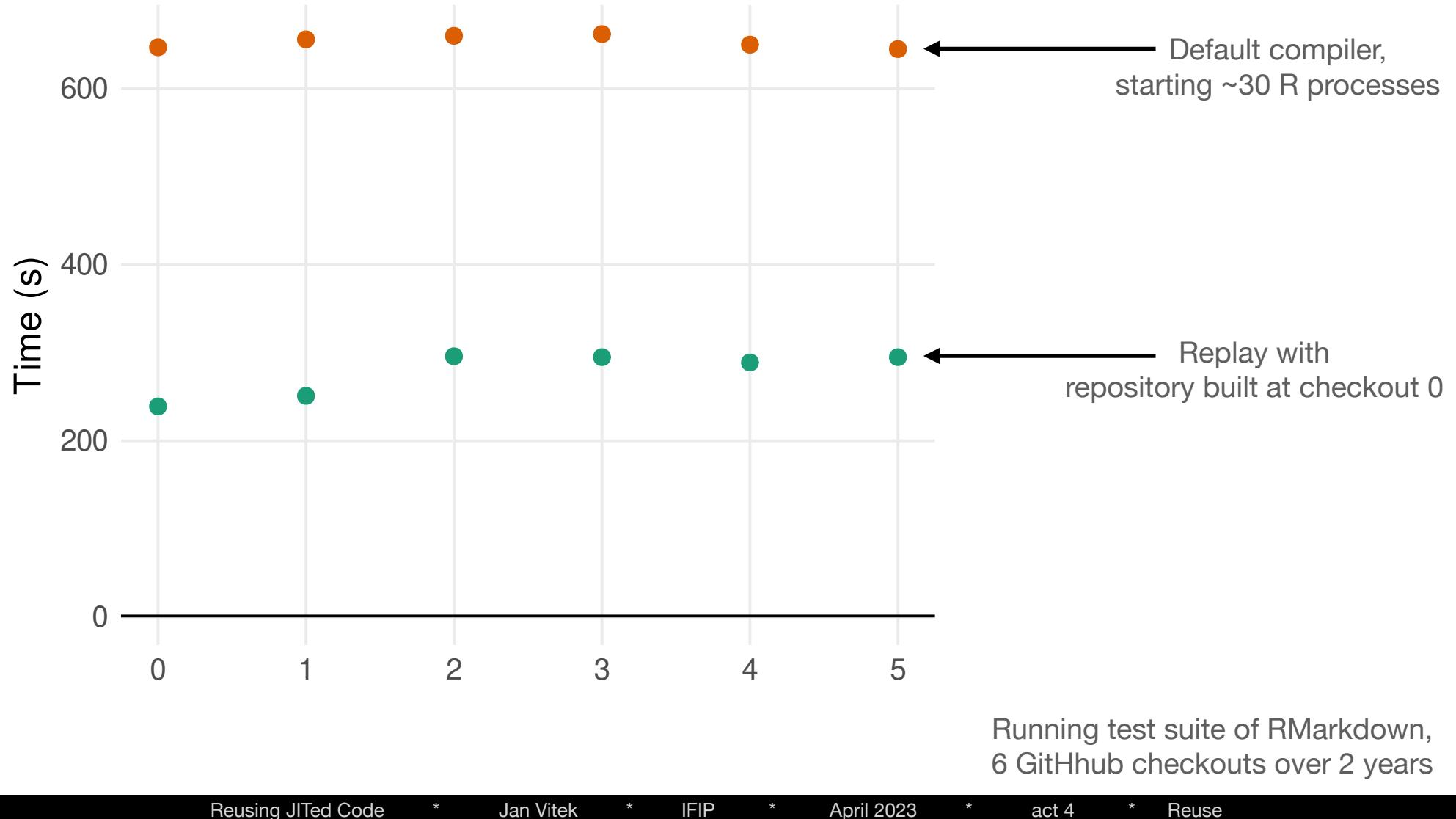
assume [r1:int] ⚓ L2

assume [r3 = 0xADD] ⚓ L2

version V3(int\str):

assume [r1:int, r2~int] ⚓ L2





Conclusions



- JITs speculate for performance
- Reusing JITed code requires reifying compiler assumptions
- Warmup greatly reduced with a JITed code repository
- Opens many research directions