# Deegen: A LLVM-based Compiler-Compiler for Dynamic Languages

Haoran Xu
haoranxu@stanford.edu
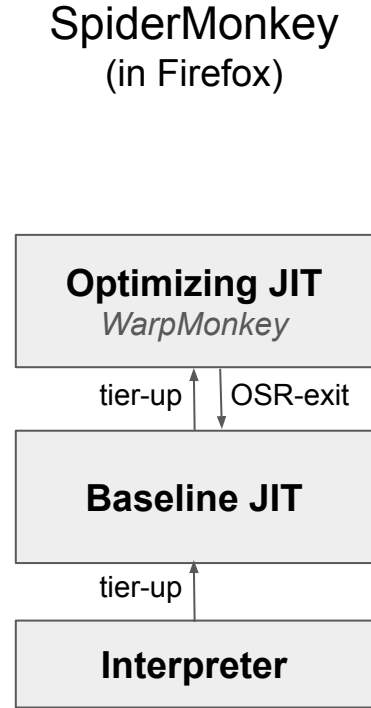
Fredrik Kjolstad
kjolstad@cs.stanford.edu

Stanford University

# Dynamic Languages

- High productivity thanks to dynamic typing.
- But also poor runtime performance on a naive VM implementation.
- And building a good VM is hard…

# Writing a good VM is hard

## JavaScriptCore
(in Safari)

| Heavyweight Opt. JIT *FTL* |
|---|

tier-up ↑     OSR exit ↓

| Lightweight Opt. JIT *DFG* |
|---|

tier-up ↑   OSR-exit ↓

| Baseline JIT |
|---|

tier-up ↑

| Interpreter *LLInt* |
|---|

## V8
(in Chrome)

| Optimizing JIT *TurboFan* |
|---|

tier-up ↑   OSR-exit ↓

| Baseline JIT *Sparkplug* |
|---|

tier-up ↑

| Interpreter *Ignition* |
|---|

## SpiderMonkey
(in Firefox)

| Optimizing JIT *WarpMonkey* |
|---|

tier-up ↑   OSR-exit ↓

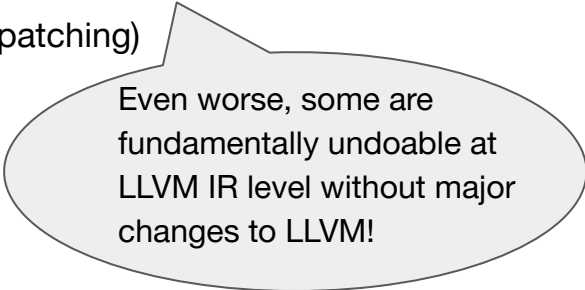| Baseline JIT |
|---|

tier-up ↑

| Interpreter |
|---|

\* OSR-exit: the process of bailing out from speculatively optimized JIT'ed code and fallback to interpreter / generic JIT'ed code, also known as deoptimization

# Can we use LLVM?

- Obviously, I'm not the first to have this idea
  - Unladen Swallow (for Python, inactive since 2010)
  - Rubinius (for Ruby, inactive since 2020)
  - LLVMLua (for Lua, inactive since 2012)
  - …
- Many attempts, but limited outreach to mainstream use
- Why?

# The Problems

- LLVM compilation is slow
  - But for a JIT, fast compilation is critical

- No direct support for the important domain-specific optimizations
  - Inline Caching / Self-Modifying Code (dynamic patching)
  - Dynamic Type Related Optimization
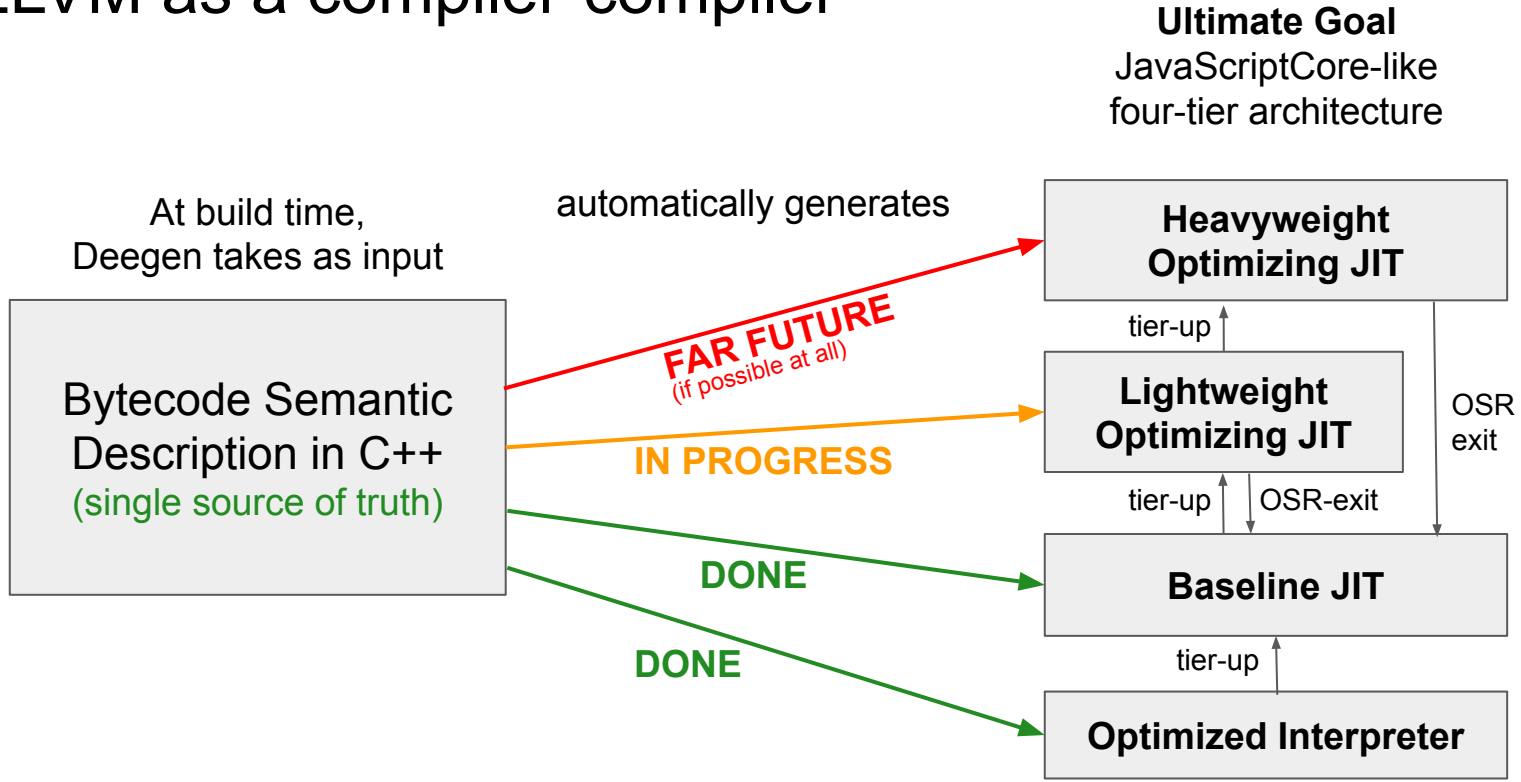  - Tiering-up / OSR-Exit
  - …

Even worse, some are fundamentally undoable at LLVM IR level without major changes to LLVM!

# Core Idea

- Do not use LLVM as a compiler
- Use LLVM as a compiler-compiler!

# LLVM as a compiler-compiler

**Ultimate Goal**
JavaScriptCore-like
four-tier architecture

At build time,
Deegen takes as input

automatically generates

Bytecode Semantic
Description in C++
(single source of truth)

**FAR FUTURE**
(if possible at all)

**IN PROGRESS**

**DONE**

**DONE**

**Heavyweight
Optimizing JIT**

tier-up

**Lightweight
Optimizing JIT**

OSR
exit

tier-up        OSR-exit

**Baseline JIT**

tier-up

**Optimized Interpreter**

# LuaJIT Remake

- Standard-compliant VM for Lua 5.1
- Bytecode execution engine generated automatically by Deegen
  - Optimized interpreter
  - Baseline JIT compiler

- VM design not identical
  - Most importantly, we have inline caching optimization (powered by Deegen)

# Performance Summary

- Interpreter-only performance
  - 31% faster than LuaJIT interpreter, 179% faster than PUC Lua
- Baseline JIT compilation cost
  - Negligible (19 million Lua bytecode/s)
- Baseline JIT performance
  - 34% slower than LuaJIT optimizing JIT, 360% faster than PUC Lua

# Bytecode Semantic Definition Example

```
1   void Add(TValue lhs, TValue rhs) {
2     if (!lhs.Is<tDouble>() || !rhs.Is<tDouble>()) {
3       ThrowError("Can't add!");
4     } else {
5       double res = lhs.As<tDouble>() + rhs.As<tDouble>();
6       Return(TValue::Create<tDouble>(res));
7     }
8   }
```

Deegen API

Defined by user, but understood by Deegen

# Bytecode Semantic Definition Example, Continued

```
1   void AddContinuation(TValue /*lhs*/, TValue /*rhs*/) {
2     Return(GetReturnValueAtOrd(0));
3   }
4   void Add(TValue lhs, TValue rhs) {
5     if (!lhs.Is<tDouble>() || !rhs.Is<tDouble>()) {
6       /* we want to call metamethod now */
7       HeapPtr<FunctionObject> mm = GetMMForAdd(lhs, rhs);
8       MakeCall(mm, lhs, rhs, AddContinuation);
9       /* MakeCall never returns */
10    } else {
11      double res = lhs.As<tDouble>() + rhs.As<tDouble>();
12      Return(TValue::Create<tDouble>(res));
13    }
14  }
```

Arbitrary runtime call, not understood by Deegen

Deegen API

Control transfers to continuation functor when call returns

# Bytecode Specification Language

```
1   DEEGEN_DEFINE_BYTECODE(Add) {
2     Operands(
3       BytecodeSlotOrConstant("lhs"),
4       BytecodeSlotOrConstant("rhs")
5     );
6     Result(BytecodeValue);
7     Implementation(Add);
8     Variant(
9       Op("lhs").IsBytecodeSlot(),
10      Op("rhs").IsBytecodeSlot()
11    );
12    Variant(
13      Op("lhs").IsConstant<tDoubleNotNaN>(),
14      Op("rhs").IsBytecodeSlot()
15    );
16    Variant(
17      Op("lhs").IsBytecodeSlot(),
18      Op("rhs").IsConstant<tDoubleNotNaN>()
19    );
20  }
```

Deegen understands the type system, and will do optimizations using this info

Also supports static quickening based on type assumption (not shown)

# User-Friendly Bytecode Builder API

```
bytecodeBuilder.CreateAdd({
    .lhs = Local(1),
    .rhs = Cst<tDouble>(123.4),
    .output = Local(2)
});
```

# Actual Disassembly of AddVV bytecode
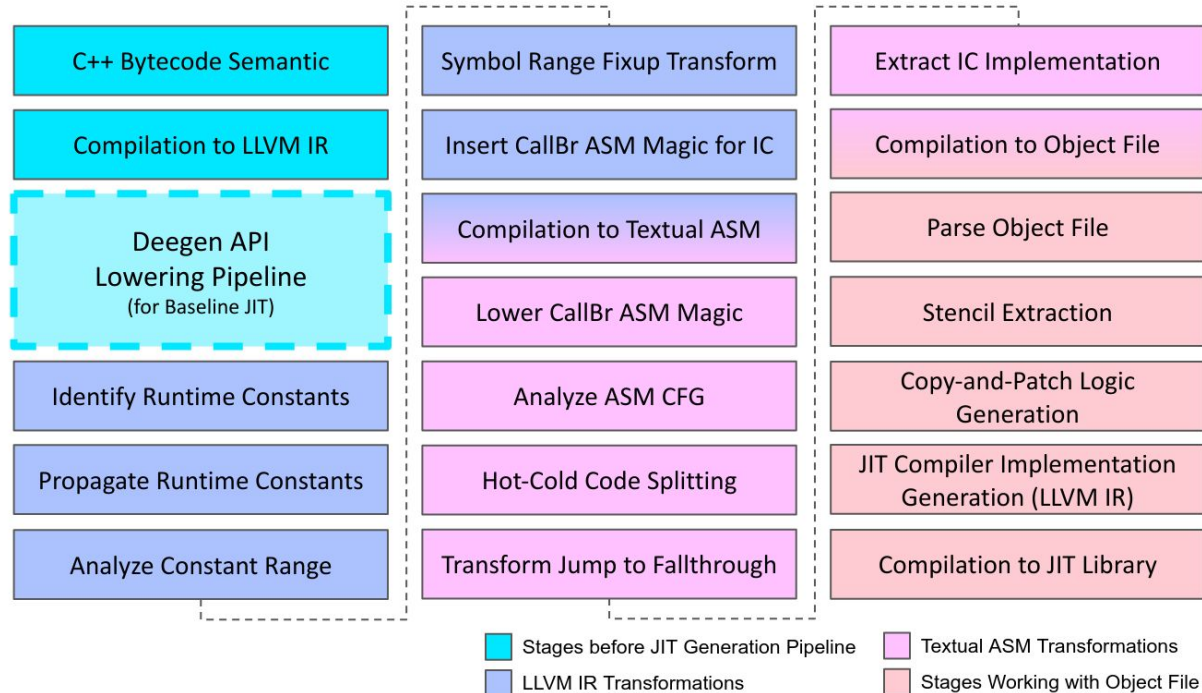
```asm
1   __deegen_interpreter_op_Add_0:
2       # decode 'lhs' from bytecode stream
3       movzwl      2(%r12), %eax
4       # decode 'rhs' from bytecode stream
5       movzwl      4(%r12), %ecx
6       # load the bytecode value at slot 'lhs'
7       movsd       (%rbp,%rax,8), %xmm1
8       # load the bytecode value at slot 'rhs'
9       movsd       (%rbp,%rcx,8), %xmm2
10      # check if either value is NaN
11      # Note that due to our boxing scheme,
12      # non-double value will exhibit as NaN when viewed as double
13      # so this checks if input has double NaN or non-double value
14      ucomisd     %xmm2, %xmm1
15      # branch if input has double NaN or non-double values
16      jp          .LBB0_1
17      # decode the destination slot from bytecode stream
18      movzwl      6(%r12), %eax
19      # execute the add
20      addsd       %xmm2, %xmm1
21      # store result to destination slot
22      movsd       %xmm1, (%rbp,%rax,8)
23      # decode next bytecode opcode
24      movzwl      8(%r12), %eax
25      # advance bytecode pointer to next bytecode
26      addq        $8, %r12
27      # load the interpreter function for next bytecode
28      movq        __deegen_interpreter_dispatch_table(,%rax,8), %rax
29      # dispatch to next bytecode
30      jmpq        *%rax
31  .LBB0_1:
32      # branch to automatically generated slowpath (omitted)
33      jmp         __deegen_interpreter_op_Add_0_quickening_slowpath
```

# The Baseline JIT Tier

- Completely free for a language implementer:
  - No additional input required.
  - Everything generated automatically from the bytecode semantics.
- Features:
  - Extremely fast compilation speed
  - Good machine code quality (under design constraints of baseline JIT)
  - Almost all optimizations used in JavaScriptCore's baseline JIT
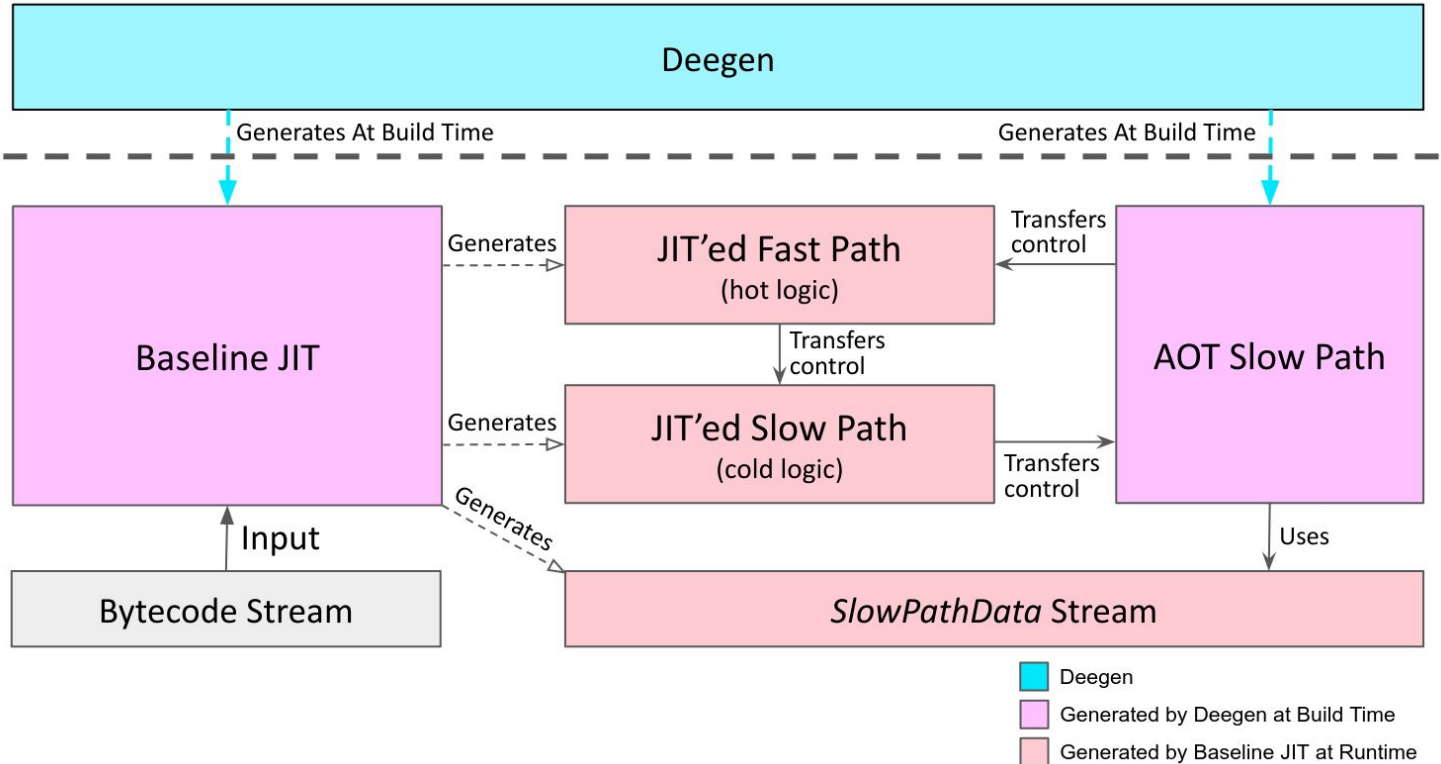
# The Baseline JIT Tier

- Generated automatically via a sophiscated build-time pipeline



| | | |
|---|---|---|
| C++ Bytecode Semantic | Symbol Range Fixup Transform | Extract IC Implementation |
| Compilation to LLVM IR | Insert CallBr ASM Magic for IC | Compilation to Object File |
| Deegen API Lowering Pipeline (for Baseline JIT) | Compilation to Textual ASM | Parse Object File |
| Identify Runtime Constants | Lower CallBr ASM Magic | Stencil Extraction |
| Propagate Runtime Constants | Analyze ASM CFG | Copy-and-Patch Logic Generation |
| Analyze Constant Range | Hot-Cold Code Splitting | JIT Compiler Implementation Generation (LLVM IR) |
| | Transform Jump to Fallthrough | Compilation to JIT Library |

Legend:
- Stages before JIT Generation Pipeline
- LLVM IR Transformations
- Textual ASM Transformations
- Stages Working with Object File

# The Baseline JIT Tier

- Use Copy-and-Patch to generate code.
- Inline Caching as the only high-level optimization
  - As it is the only high-level optimization that can be performed without sacrificing startup delay
- However, many low-level optimizations
  - Runtime-constant propagation (aka, binding-time analysis)
  - Self-modifying-code-based IC implementation for best perf
  - Inline Slab optimization for IC
  - Hot-cold splitting
  - Tail-jump elimination
  - …

# Baseline JIT Architecture (except Inline Caching)

# Example: generated code for Add

```
fast_path:
   0: f2 0f 10 8d ** ** ** **    movsd    $ 1 (%rbp), %xmm1
   8: f2 0f 10 95 ** ** ** **    movsd    $ 2 (%rbp), %xmm2
  10: 66 0f 2e ca                ucomisd  %xmm2, %xmm1              1  lhsSlot * 8
  14: 0f 8a ** ** ** **          jp        3                       2  rhsSlot * 8
  1a: f2 0f 58 ca                addsd    %xmm2, %xmm1             3  slow_path
  1e: f2 0f 11 8d ** ** ** **    movsd    %xmm1, $ 4 (%rbp)        4  outputSlot * 8
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -  5  slowPathDataOffset
slow_path:                                                         6  __Add_slowpath
   0: 41 bc ** ** ** **          movl     $ 5 , %r12d
   6: 4c 03 63 30                addq     0x30(%rbx), %r12
   a: e9 ** ** ** **             jmp       6
```

# Closure Thoughts

- What is Deegen's #1 contribution?
  - Research novelty? Definitely a contribution, but not #1 IMO…
- What is LLVM's #1 contribution to the world?
  - The engineering that puts together decades of compiler research into a reusable infrastructure for static languages
- … that's also the story I dream for Deegen …
  - The engineering that recollects the $$$$ lessons of JSC, V8, … into a reusable infrastructure for dynamic languages
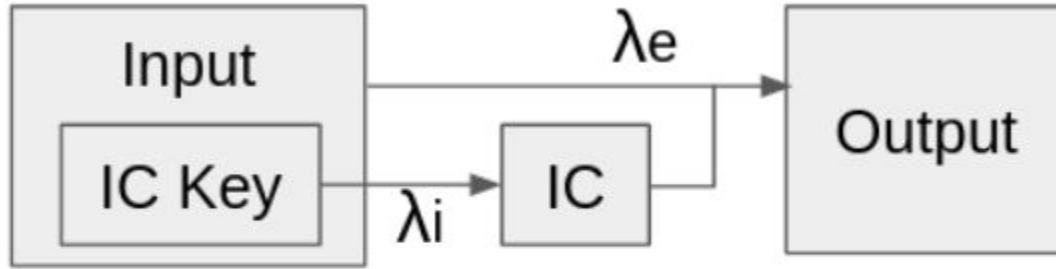  - Very hard, still very far away, but we are at a good start :)

# Extra Slides

# Inline Caching

- "The most important optimization" —JavaScriptCore dev
- Key observation: certain values can be well-predicted
  - For code `f()`, "`f`" likely holds the same function
  - Many objects are used like C structs, so a property access site (e.g., "`employee.name`") likely to see objects with the same "structure".
- Cache the seen value and computation result at use site ("inline" caching)
- If next time we see the same value, can skip redundant computation
  - For call, can skip the check that the object is indeed a function, and the load of the code pointer from the function
  - For object property access, combined with hidden class, can skip the hash table lookup and directly know where the property is

# Inline Caching in Deegen

- Deegen understands calls, but not objects
  - Object semantics drastically differ per language
  - Impossible to provide a generic and ideal implementation
  - So should not be hardcoded by Deegen
- Call inline caching
  - Automatic in Deegen, no user intervention
- Object property inline caching
  - Achieved by Generic Inline Caching API
  - Requires user to use the API to express IC semantics

# Generic Inine Caching API



$\lambda i$ : expensive but idempotent computation

$\lambda e$: cheap computation based on the input and the result of the idempotent step

Computation eligible for inline caching can be characterized as above.

# Generic Inine Caching API

- Idea: use C++ lambda to represent computation
- Body lambda
    - Represents the overall computation
- Effect lambda
    - Defined inside the body lambda, can have multiple
    - Represents an effectful computation
- That is, all computation in the body lambda must be idempotent. Effectful computation must be done within an effect lambda.

# Inline Caching Example: TableGetById

- TableGetById
- Get a fixed string property from the table
- **e.g.,** `employee.name`, `animal.weight`
- One of the most common operations on object.

```
1   void TableGetById(TValue tab, TValue key) {
2       // Let's assume 'tab' is indeed a table for simplicity.
3       HeapPtr<TableObject> t = tab.As<tTable>();
4       // And we know 'key' must be string since the index of
5       // TableGetById is required to be a constant string
6       HeapPtr<String> k = key.As<tString>();
7       // Call API to create an inline cache
8       ICHandler* ic = MakeInlineCache();
9       HiddenClassPtr hc = t.m_hiddenClass;
10      // Make the IC cache on key 'hc'
11      ic->Key(hc);
12      // Specify the IC body (the function 'λ')
13      Return(ic->Body([=] {
14          // Query hidden class to get value slot in the table
15          // This step is idempotent due to the design of hidden class
16          int32_t slot = hc->Query(k);
17          // Specify the effectful step (the function 'λ_e')
18          if (slot == -1) {       // not found
19              return ic->Effect([] { return NilValue(); }
20          } else {
21              return ic->Effect([=] { return t->storage[slot]; });
22          }
23      });
24  }
```

**The Body Lambda** → (line 13: `ic->Body([=]`)

**Two Effect Lambdas** → (lines 19 and 21: `ic->Effect`)

**Value defined in body lambda Treated as result from idempotent computation** → (line 21: `slot`)

**Value defined outside, sees fresh value every time** → (line 21: `t->storage`)

# TableGetById: Interpreter Logic Disassembly

```
__deegen_interpreter_op_TableGetById_0_fused_ic_3:
    pushq    %rax
    movzwl   2(%r12), %eax                    # decode the src slot from bytecode
    movq     (%rbp,%rax,8), %r9               # load the src TValue from stack
    cmpq     %r15, %r9                        # check if it is a heap entity
    jbe      .LBB5_9                          # if not, branch to slow path (omitted)
    movzwl   6(%r12), %r10d                   # Decode the dst slot from bytecode
    movl     8(%r12), %edi
    addq     %rbx, %rdi                       # Get metadata struct (holding the inline cache for this bytecode)
    movl     %gs:(%r9), %ecx                  # Load hidden class (safe as we have checked it's a heap entity)
    cmpl     %ecx, (%rdi)                     # Check if inline cache hits
    jne      .LBB5_5                          # If not, branch to slow path (omitted)
    movslq   5(%rdi), %rax                    # IC directly tells us the slot holding the property in the object
    movq     %gs:16(%r9,%rax,8), %rax         # Load that slot in the object
    movq     %rax, (%rbp,%r10,8)              # Store the result back to dst slot in the stack frame
    movzwl   12(%r12), %eax                   # Dispatch to next bytecode
    addq     $12, %r12
    movq     __deegen_interpreter_dispatch_table(,%rax,8), %rax
    popq     %rcx
    jmpq     *%rax
```

# Baseline JIT Inline Caching Design

# Further Reading

- My Blog:
  - sillycross.github.io
- Blog post titles:
  - Building the fastest Lua interpreter automatically
  - Building a baseline JIT for Lua automatically
- LuaJIT Remake Github repo:
  - https://github.com/luajit-remake/luajit-remake