

Native Law Judgements

Benjamin Summers

February 19, 2026

Abstract

This document describes the design and implementation of the default function evaluation engine of the native PLAN runtime system.

1 Introduction

In order to achieve a system which truly has zero dependencies, our PLAN runtime system is written in assembly. This makes the code quite expensive to maintain. Worse, the runtime system is designed to have up-times in the decades, which makes it impossible to do a live upgrade of the runtime system itself. This constraint adds a number of difficult and unconventional engineering challenges.

Our high level approach to these challenges is to keep the runtime system as small as possible, and to instead move as much code into XPLAN, which can be written in a much higher-level language and which can be upgraded dynamically.

One of the most significant subsystems which we move into XPLAN is the actual optimization logic used to run functions. However, a number of non-trivial tasks must still be possible using only the basic evaluation system:

1. We need to be able to run a significant body of code in order to test the runtime system itself.
2. The optimization engine will be written in source code in a high-level language. We will need to be able to load a compiler and use it to build that optimization engine.
3. The compiler for this language will also be written in a high level language, so we will need to be able to run it against itself for bootstrapping.
4. In order to produce the actual seed-file that we need for bootstrapping, we need to be able to serialize the bootstrapping compiler.

This tension between the need to minimize mechanical complexity, and the need to achieve a usable performance baseline results in a somewhat unique implementation strategy, which we will document in detail in the rest of this document.

2 Lazy Graph Reduction

PLAN is a lazy graph reduction engine built using super-combinators. A system like this, in theory, doesn't do any actual computation when a function is called. Functions just work something like a template expansion, and the graph reduction engine does all of the actual work.

Here's a concrete example:

```
foo :: Bool -> Int -> Int -> Int
foo x y z = if x then y+z else z+z
```

If `foo` is invoked as `(foo True 3 4)`, the actual result of running the function is not 6, but tree of thunks, which is then processed by the graph reduction engine in order to find a result:

```
((((foo True) 3) 4)
((foo[True] 3) 4)
(foo[True] 3) 4)
(((If True) ((Add 3) 3)) ((Add 4) 4)) -- expansion of foo
((If[True] ((Add 3) 3)) ((Add 4) 4))
(If[True] ((Add 3) 3)) ((Add 4) 4))
((Add 3) 3) -- expansion of If
(Add[3] 3)
6 -- expansion of Add
```

Producing and consuming all of these thunks is highly inefficient, so sophisticated runtime systems figure out which subset of the work will always be performed, and just do all of that eagerly. But this requires somewhat sophisticated analysis and optimization work.

In order to keep the native runtime system simple, we avoid all of this work and instead adopt the naive template-expansion approach. But, in order to produce something with a usable base-line of performance, we also perform a number of rote optimizations which only require trivial analysis.

3 Optimizing Template Expansion

Here's our example again:

```
= (mapMaybe f o)
| If (Nil o) 0
| o (f (Ix 0 o))
```

Ignoring the issues around needing to make everything legible to GC, the most stupid possible approach to template expansion would look something like this:

```
Obj mapMaybe(Obj f, Obj o) {
    Obj result =
        App(App(App(If, App(Nil,o),
                    o),
                  App(o,
                      App(f,
                          App(App(Ix,o),
                              o))));
    tailcall tmp.execute()
}
```

This approach requires a very large number of small thunks to be allocated, and each of those thunks must be separately processed by the graph reduction engine.

To improve on this, we use a series of improvements on this basic output in order to get much better performance.

3.1 Optimization 1: Bigger Thunks

We can do the same thing with fewer allocations, and make life a bit easier for the graph reduction engine by constructing bigger thunks instead of only use Apps.

```
Obj mapMaybe(Obj f, Obj o) {
    Obj result =
        App3(If, App(Nil,o),
              o,
              App(o, App(f, App2(Ix, 0, o))));
    tailcall tmp.execute()
}
```

3.2 Optimization 2: Direct Closure Construction

Notice that the `(o (f (Ix 0 o)))` expression is always undersaturated, like an expression like `(Add 3)` would be. So the result will always be a closure.

Producing a thunk in this case results in a lot of wasted work. The thunk will examine the head, determine that the call is undersaturated, allocate a closure, and update the thunk to point to the closure.

As long as the function call is to a known constant, it is trivial to determine that the call is undersaturated, which allows us to skip constructing this thunk, and instead directly construct a closure.

```
Obj mapMaybe(Obj f, Obj o) {
    Obj result =
        App3(If, App(Nil,o),
              o,
              C1z1(o, App(f, App2(Ix, 0, o))));
    tailcall tmp.execute()
}
```

3.3 Optimization 3: Single Allocation

Because the template expansion for each function is always the exact same shape, we can avoid doing many small allocations, and instead just do one big allocation.

Not only does this save work by avoiding repeated allocation, but it also means that the work of actually filling in the template does not have to make sure to make all of its pointers available to the garbage collector. No GC will happen during this initialization process.

```
Obj mapMaybe(Obj f, Obj o) {
    save(f, o);
    void *buf = alloc(22);
    restore(f, o);

    memcpy(template, buf);

    buf[2] = Nil; // from constants array
    buf[3] = o; // from register
    buf[5] = Ix; // from constants array
    buf[7] = o; // from register
    buf[10] = f; // from register
    buf[11] += buf; // convert offset to pointer
    buf[15] += buf; // convert offset to pointer
    buf[17] = If; // from constants array
    buf[18] += buf; // a: convert offset to pointer
    buf[20] += buf; // d: convert offset to pointer

    Obj *result = (Obj*) &result[17]; // 17 b/c first word is GC header
    tailcall result.execute()
}
```

3.4 Optimization 4: Precomputed Template

Not only is the size of the resulting graph constant, but so is its structure. If we pre-compute the shape of the template, we can simply copy the whole thing to our reserved space on the heap.

There are only a few changes that need to be made to this template after it is copied:

- Pointers between nodes internal to the graph need to be updated to point to the new node locations.

- Reference to variables must be filled in.

- Because references to constants (except direct nats) can be moved around by GC, we need to maintain a table of constants, and fill each of these in by replacing them with a reference into this table.

```
Obj mapMaybe(Obj f, Obj o) {
    save(f, o);
    void *buf = alloc(22);
    restore(f, o);

    memcpy(template, buf);

    buf[2] = Nil; // from constants array
    buf[3] = o; // from register
    buf[5] = Ix; // from constants array
    buf[7] = o; // from register
    buf[10] = f; // from register
    buf[11] += buf; // convert offset to pointer
    buf[15] += buf; // convert offset to pointer
    buf[17] = If; // from constants array
    buf[18] += buf; // a: convert offset to pointer
    buf[20] += buf; // d: convert offset to pointer

    Obj *result = (Obj*) &result[17]; // 17 b/c first word is GC header
    tailcall result.execute()
}
```

3.5 Optimization 5: Code Generation

As you can see from the previous example, the actual template-filling logic is just a straight-line sequence of simple memory operations, and none of this needs to be interleaved with garbage collection.

As a result, it is quite simple to just generate this machine code for each law at law construction time.

In particular, we only need to be able to construct 4 assembly expressions. If constants=r11 and output=r12, then:

```
1. Load a constant: mov r10, [r12+o]
2. Load from stack: mov r10, [r15+o]
3. Write a register: mov [r14+o], reg, (where reg can be: rax, r10, rdi, rsi, rdx, rcx, r8, or r9).
4. Write a argument: mov [r14+o], rdi
5. Hydrate an internal: add [r14+o], r14
```

If we use 32-bit offsets for each output instead of trying to minimize the size, we can use a simple 7-byte output for each instruction, which allows a simple table to be used instead of complex encoding logic.

```
Registers:
    r15      = stack
    r14, r13 = reserved (heap registers)
    r12      = temp
    r11      = constants
    rax      = self reference
    rdi, rsi..,r12 = arguments
```

Patterns:

```
    mov r12, [r11+x] ; 4D8BA3xxxxxxxx
    mov r12, [r15+x] ; 4D8BA7xxxxxxxx
    mov [r14+x], rax ; 498986xxxxxxxx
    mov [r14+x], rdi ; 49898BExxxxxxx
    mov [r14+x], rsi ; 498986xxxxxxxx
    mov [r14+x], rdx ; 498996xxxxxxxx
    mov [r14+x], r8 ; 4D8986xxxxxxxx
    mov [r14+x], r9 ; 4D8998xxxxxxxx
    mov [r14+x], r10 ; 4D8996xxxxxxxx
    mov [r14+x], r11 ; 4D8998xxxxxxxx
    mov [r14+x], r12 ; 4D89A6xxxxxxxx
    add [r14+x], r14 ; 4D01B6xxxxxxxx
```

```
    Obj *result = (Obj*) &result[17]; // 17 b/c first word is GC header
    tailcall result.execute()
```

3.6 Optimization 6: Avoid Thunk Update

Loops in function programming languages are encoded using tail recursion, which makes tail recursion optimization an essential feature of any purely functional programming language like PLAN.

Tail call optimization is a bit different in a minimalist, purely functional language like PLAN, because things like `if` and `case` and just normal functions.

Lazy evaluation essentially works by repeatedly simplifying the outer expression until it is "done" and then repeating the process for the inner expressions. This iterative process automatically provides constant-space tail recursion except for one tricky-edge case: thunk caching.

After each thunk is evaluated, it is modified so that subsequent references return the pre-computed value, instead of running the computation again. But, in order to do this modification, we need to run a piece of logic after the evaluation a stack frame, and breaks tail recursion.

The problem is that we end up with a long chain of tiny stack frame which just update thunks. If we could find some way to use an alternate thunk structure which doesn't perform the update, all of these tiny stack frames would go away.

Consider the following example:

```
= (sillyAdd a b)
| Seq a
| Ifz b a
| sillyAdd Inc-a Dec-b
```

```
= (sillyDouble a)
| sillyAdd a a
```

When sillyAdd is run, we know that the call to sillyAdd will always be evaluated immediately, so the thunk will definitely only be evaluated once. In this case, we can trivially use a non-updating thunk and recover tail recursion.

In 'sillyAdd', that means that (`Seq a ..`) doesn't need to update, but what about (`Ifz b a ..`)? Well, we know that Seq evaluates each argument exactly once, so any expression passed as an argument to `(Add 3)` would be a thunk update. And the same with `'Ifz'`. Ifz will evaluate each of its arguments at most 1 time, so there is no need for the thunk update.

This chain of reasoning eliminates the thunk updates for the whole chain of calls leading up the the recursion, and the result is that we recover tail call optimization.

Doing this analysis in general requires a significant amount of work. However, in XPLAN, things like Seq and Ifz are treated as primitives, so we can simply hard-code this information on functions which are simple wrappers for primops (see the later section on recognizing primop wrappers).

Once this information is available, the process is trivial. Whenever we see a saturated call to a function, remember which arguments are evaluated at most once. When generating a thunk for such a function, use the non-updating version.

4 Template Layout

First of all, we special case trivial functions which just return an argument or a constant value. These functions don't require any allocations at all.

Computing the template is fairly easy, we just need to figure out what all of the nodes are and where they belong in the template. There are only a few cases:

1. A saturated or unknown function call will be a thunk, and will require args+3 words.
2. An undersaturated function call will be a closure, and that requires args+2 words.
3. A constants value or a variable reference requires no space in the template, just a reference from the containing expression.

Because the template expansion for each function is always the exact same shape, we can avoid doing many small allocations, and instead just do one big allocation.

Not only does this save work by avoiding repeated allocation, but it also means that the work of actually filling in the template does not have to make sure to make all of its pointers available to the garbage collector. No GC will happen during this initialization process.

```
Obj mapMaybe(Obj f, Obj o) {
    save(f, o);
    void *buf = alloc(22);
    restore(f, o);

    memcpy(template, buf);

    buf[2] = Nil; // from constants array
    buf[3] = o; // from register
    buf[5] = Ix; // from constants array
    buf[7] = o; // from register
    buf[10] = f; // from register
    buf[11] += buf; // convert offset to pointer
    buf[15] += buf; // convert offset to pointer
    buf[17] = If; // from constants array
    buf[18] += buf; // a: convert offset to pointer
    buf[20] += buf; // d: convert offset to pointer

    Obj *result = (Obj*) &result[17]; // 17 b/c first word is GC header
    tailcall result.execute()
}
```

3.4 Optimization 4: Precomputed Template

Not only is the size of the resulting graph constant, but so is its structure. If we pre-compute the shape of the template, we can simply copy the whole thing to our reserved space on the heap.

There are only a few changes that need to be made to this template after it is copied:

- Pointers between nodes internal to the graph need to be updated to point to the new node locations.

- Reference to variables must be filled in.

- Because references to constants (except direct nats) can be moved around by GC, we need to maintain a table of constants, and fill each of these in by replacing them with a reference into this table.

```
Obj mapMaybe(Obj f, Obj o) {
    save(f, o);
    void *buf = alloc(22);
    restore(f, o);

    memcpy(template, buf);

    buf[2] = Nil; // from constants array
    buf[3] = o; // from register
    buf[5] = Ix; // from constants array
    buf[7] = o; // from register
    buf[10] = f; // from register
    buf[11] += buf; // convert offset to pointer
    buf[15] += buf; // convert offset to pointer
    buf[17] = If; // from constants array
    buf[18] += buf; // a: convert offset to pointer
    buf[20] += buf; // d: convert offset to pointer

    Obj *result = (Obj*) &result[17]; // 17 b/c first word is GC header
    tailcall result.execute()
}
```

3.2 Optimization 2: Direct Closure Construction

Notice that the `(o (f (Ix 0 o)))` expression is always undersaturated, like an expression like `(Add 3)` would be. So the result will always be a closure.

Producing a thunk in this case results in a lot of wasted work. The thunk will examine the head, determine that the call is undersaturated, allocate a closure, and update the thunk to point to the closure.

As long as the function call is to a known constant, it is trivial to determine that the call is undersaturated, which allows us to skip constructing this thunk, and instead directly construct a closure.

```
Obj mapMaybe(Obj f, Obj o) {
    Obj result =
        App3(If, App(Nil,o),
              o,
              App(o, App(f, App2(Ix, 0, o))));
    tailcall tmp.execute()
}
```

3.3 Optimization 3: Single Allocation

Because the template expansion for each function is always the exact same shape, we can avoid doing many small allocations, and instead just do one big allocation.

Not only does this save work by avoiding repeated allocation, but it also means that the work of actually filling in the template does not have to make sure to make all of its pointers available to the garbage collector. No GC will happen during this initialization process.

```
Obj mapMaybe(Obj f, Obj o) {
    save(f, o);
    void *buf = alloc(22);
    restore(f, o);

    memcpy(template, buf);

    buf[2] = Nil; // from constants array
    buf[3] = o; // from register
    buf[5] = Ix; // from constants array
    buf[7] = o; // from register
    buf[10] = f; // from register
    buf[11] += buf; // convert offset to pointer
    buf[15] += buf; // convert offset to pointer
    buf[17] = If; // from constants array
    buf[18] += buf; // a: convert offset to pointer
    buf[20] += buf; // d: convert offset to pointer

    Obj *result = (Obj*) &result[17]; // 17 b/c first word is GC header
    tailcall result.execute()
}
```

3.4 Optimization 4: Precomputed Template

Not only is the size of the resulting graph constant, but so is its structure. If we pre-compute the shape of the template, we can simply copy the whole thing to our reserved space on the heap.

There are only a few changes that need to be made to this template after it is copied:

- Pointers between nodes internal to the graph need to be updated to point to the new node locations.

- Reference to variables must be filled in.

- Because references to constants (except direct nats) can be moved around by GC, we need to maintain a table of constants, and fill each of these in by replacing them with a reference into this table.

7 Recording Multiplicity Information

We can encode multiplicity information as a simple bit-mask. The lowest bit being set indicates that the first argument is unshared, etc.

However, where in the law allocation box should this information be stored?

8 Storing the Code

We can construct a NAT which stores the execution logic for a law, and simply store it in one of the metadata slots for a LAW.

However, we do need to allocate once in each routine, and we don't want the code to move around from underneath us, so we need some way to store this outside of the GC heap.

Fortunately, I have been exploring ideas for moving large nats off-heap anyways.

I believe that this can be achieved by adding a mark bit to the GC header, and adding a bit to the NAT tag which indicates whether or not it lives in a finalized pin or not (second or third GC generation).

Each thread will need to keep a table of large NATs (probably as a backwards linked list stored in normal GC memory). Whenever we see a reference to a local bignum, we must mark it. And, at the end of GC, we need to traverse the list of known bignums and free each one which has not been marked.

Erlang, like many other systems, does something like this, because otherwise your big nats need to be copied around on every GC.

This system would also mean that the code that we generate for laws would not move around during GC, which solves a lot of problems.

TODO: should the actual code-pointer point directly into this NAT. Is that safe?

9 Storing the Constants Table

The code will also need access to a table of constants, and this will need to be placed in the law somewhere as well.

I suppose we can do this all with a single metadata slot:

```
0[multiplicity template-and-code constant1 constant2..]
```

Where the code NAT directly includes the output template as a constant.