

Compiling for Runtime Code Generation

Frederick Smith, Dan Grossman, Greg Morrisett, Luke Hornof, and Trevor Jim

Cornell University, {fms,danieljg,jgm}@cs.cornell.edu

Transmeta Corporation, hornof@transmeta.com

AT&T Research, trevor@research.att.com

Abstract. Cyclone is a programming language that provides explicit support for dynamic specialization based on runtime code generation. To generate specialized code quickly, our Cyclone compiler uses a template based strategy in which pre-compiled code fragments are stitched together at runtime. To achieve good performance, the pre-compiled fragments must be optimized. This paper describes a principled approach to achieving such optimizations. In particular, we generalize standard flow-graph intermediate representations to support templates, define a formal mapping from (a subset of) Cyclone to this representation, and describe a data-flow analysis framework that supports standard optimizations.

1 Introduction

Runtime code generation (RTCG) is a kind of program specialization that optimizes for short-lived invariants by dynamically generating specialized code. In this setting, where invariants may last only seconds or minutes, compilation time can easily dominate any benefit. Therefore, rapid compilation is critical. *Template-based RTCG* [3, 13, 2] generates code by combining copies of pre-compiled code fragments (templates). Compared to other approaches (see Section 7), template-based RTCG enjoys faster compilation but less optimization.

In earlier work [8], we developed *Cyclone*, a type-safe C-like language with support for explicit RTCG. In that work, we proved the soundness of Cyclone's type system and presented a simple stack-based compiler that translates Cyclone to TAL/T, an extension of Typed Assembly Language [11]. The type-checker for TAL/T statically guarantees that, in addition to standard type-safety properties, no type-unsafe code will be generated at run time. In the context of reliable or secure systems (*e.g.*, operating systems), this guarantee is crucial for assurance.

Though Cyclone allows programmers to dynamically generate specialized code, our previous compiler performed almost no optimization. For instance, we placed all variables on the stack instead of performing register allocation. As a result, the performance of the dynamically generated code was relatively poor compared to statically generated, but highly-optimized, code. In contrast, the Tempo template-based RTCG system [3, 13] is able to achieve high performance by generating C code with a clever use of labels (to define template boundaries), and then passing the resulting code to GCC. Though an expedient approach, the Tempo designers had to trick GCC into compiling template-based code, and as

a result, the compiler does not always generate correct code. Furthermore, the output of GCC cannot be (easily) type-checked. Consequently, Tempo cannot be used in settings where reliability or security are just as important as performance.

We therefore set out to adapt traditional flow-graph intermediate representations, data-flow analyses (such as liveness analysis), and optimizations (such as register allocation) to the template-based RTCG setting, so that we might achieve both high-performance and high-assurance. This project requires a template-aware intermediate representation, a translation from Cyclone to the less structured intermediate form, and a sound and effective framework for performing data-flow-based optimization. We quickly encountered a number of subtleties that suggested these tasks were more difficult than they first appeared.

Our contributions include (1) a formal translation from a source language to an intermediate form, (2) an analysis for computing control-flow graphs for functions containing templates, and (3) a discussion about performing optimizations in the presence of RTCG. Although we describe small theoretical languages, they capture the most interesting aspects of our new optimizing RTCG compiler.

2 Overview

Cyclone is a type-safe C-like language supporting polymorphism, exceptions, and, of course, RTCG. There are four constructs that support RTCG: `codegen`, `cut`, `splice`, and `fill`. These constructs are best conceptualized as manipulating “strings” to construct the text of a new function. For example, the function `dot_gen` below takes a vector `u` and returns a function specialized for taking dot products of `u`. The boxes are used to highlight code fragments that are dynamically assembled to produce the specialized function.

```
(int[] -> int) dot_gen(int u[]) {
  return codegen (
    int dot(int v[]) {      int result = 0;
      cut { for(int i=0; i<size(u); i++)
        splice{ result += fill(u[i]) * v[fill(i)]; }; }
      return result;    });
}
```

The expression `codegen (int(dot(intv[])...)` indicates runtime code generation of a function, `dot`. The `cut {...}` suspends code generation and drops back into running code for `dot_gen`. On each loop iteration, `splice {...}` appends a copy of its contents to the body of `dot`. Furthermore, each copy is specialized using the `fill` construct to place the values of `u[i]` and `i` in the generated code. Thus, executing `dot_gen({12,24,32})` returns a function like this one:

```
int dot(int v[]) {
  int result = 0;
  result += 12 * v[0];
```

```

    result += 24 * v[1];
    result += 32 * v[2];
    return result;
}

```

In this example, RTCG allows us to eliminate the looping overhead and to fold constants into the instructions. A more aggressive version of `dot_gen` could generate different code when an element of `u` is zero, one, or a power of two.

The RTCG primitives could be implemented using string manipulation and calls to a conventional compiler. Doing so would ensure that the resulting code was well optimized. For instance, we could reasonably expect that the `result` variable would be allocated a register. However, the overhead of a conventional compiler makes it difficult to take advantage of short-lived invariants.

We therefore use a different approach, leveraging the fact that Cyclone does not permit arbitrary operations on code fragments. In particular, we pre-compile code fragments into templates—binary code fragments with “holes” for run-time values—when we initially compile the program. These templates can be quickly copied, and their holes can be filled, to produce a specialized function with low overhead. For example, the compiled `dot_gen` would use three templates: (1) for the function prologue and for initializing `result`, (2) for adding to `result` some constant value (`u[i]`) multiplied by a vector element at some constant offset (`i`), and (3) for returning `result` and the function epilogue.

To achieve good code, we must perform optimizations on the templates, such as register allocation, at compile time. But how can we perform the necessary data-flow analyses, such as live-range analysis, when templates can be dynamically assembled at run-time? The approach we suggest here is to perform data-flow analysis on the *generating* function to construct a *generating graph*—an approximation of the control-flow graph—for the functions that may be generated at run-time. Then we can use the generating graph to analyze the possible data-flow paths of the templates and optimize them.

For example, by analyzing `dot_gen`, we know that, no matter what specialized function it produces, the control flow graph will have a restricted form: It will start with one copy of the first template, followed by some number of copies of the second template (appropriately specialized with constants), and then one copy of the third. Such an analysis would allow us to compute accurate live ranges for the variables of the function, and to allocate them to registers.

In what follows, we make these informal ideas precise by formally defining a core subset of Cyclone (Section 3), an intermediate representation with explicit support for templates (Section 4), and a compiler that maps Mini-Cyclone to the intermediate representation (Section 5). We then discuss how analyses and optimizations may be applied to the intermediate representation (Section 6).

3 Mini-Cyclone

Figure 1 presents the abstract syntax for Mini-Cyclone, a subset of Cyclone that serves as this paper’s source language. The primitive values are integers (*i*)

(integers)	$i \in \mathcal{Z}$
(function names)	$f \in \text{CYCVAR}$
(variables)	$x \in \text{CYCVAR}$
(expressions)	$e ::= i \mid f \mid x \mid e_1 + e_2 \mid e(e_1, \dots, e_n) \mid \text{codegen } D \mid \text{fill } e$
(statements)	$s ::= x := e \mid s_1; s_2 \mid \text{if}(e) s_1 \text{ else } s_2 \mid \text{while}(e) s \mid \text{return } e$ $\quad \quad \quad \mid \text{cut } s \mid \text{splice } s$
(functions)	$D ::= f(x_1, \dots, x_n) s$
(programs)	$P ::= D_1, \dots, D_n$

Fig. 1. Mini-Cyclone Syntax

and function names (f). The other standard expressions are variables, addition, and function application. Similar to Cyclone, there are two expression forms for RTCG: `codegen` d generates a function and returns a pointer to it. `fill` e uses the evaluation of e as a constant in the enclosing code fragment.

The standard statements (assignment, composition, conditional, loop, and return) all have their usual meaning. `cut` s must occur within a code fragment for a generated function. Its meaning is to execute s in the generating function while the generated function is being generated. `splice` s must occur within a `cut`. It treats s as a code fragment belonging to the function that was being generated when the enclosing `cut` began executing. A copy of s is placed after the code that has already been generated.

A Mini-Cyclone function takes a fixed number of arguments and executes a body, which is just a statement. For simplicity here, we ignore types and assume that code is well-formed.

4 Target Language

Our compiler’s intermediate representation is a standard low-level block-based language with special RTCG instructions and the type information needed to generate certified code. In this paper, we focus on compiling for RTCG, so we ignore typing issues and use a target language called CIR that is sufficient for compiling Mini-Cyclone.

Figure 2 shows the syntax for CIR. A CIR program is a collection of functions. A function has a list of parameters and a list of blocks. A block has a label (l), a template name (t), an instruction sequence (ι_1, \dots, ι_n), and a control transfer (χ). The first block is the function’s entry point; its label serves as the function’s name. Assignment ($x := v$), addition ($x := v_1 + v_2$), function call ($x := v(x_1, \dots, x_n)$), return (`retn` x), and jump (`jmp` d) are standard. The conditional transfer (`jnz` $x ? d_1 : d_2$) jumps to d_1 if x is not zero, else it jumps to d_2 . The destination \circ is a “fall through” — in block B_i of a function, \circ is equivalent to the label of B_{i+1} .

A CIR function represents either code that is directly executable or a (template) container. Executable functions must not contain holes, whereas contain-

	(instruction) $\iota ::= x := v$
$i \in \mathcal{Z}$ $x \in \text{VAR}$ $l \in \text{LABEL}$ $h \in \text{HOLE}$ $t \in \text{TEMPLATE}$ $r \in \text{REGION}$ $p \in \text{INSTANCE}$	$x := v_1 + v_2$ $x := v(x_1, \dots, x_n)$ $x := [h]$ $r := \underline{\text{start}}\ l$ $p := \underline{\text{copy}}\ t\ \underline{\text{into}}\ r$ $\underline{\text{fill}}\ p.[h]\ \underline{\text{with}}\ v$ $\underline{\text{fill}}\ p_1.[h]\ \underline{\text{with}}\ p_2.l$ $x := \underline{\text{end}}\ r$
(value) $v ::= x \mid l \mid i$	(transfer) $\chi ::= \underline{\text{jmp}}\ d \mid \underline{\text{jnz}}\ x?d_1 : d_2 \mid \underline{\text{retn}}\ x$
(destination) $d ::= l \mid \circ \mid [h]$	(block) $B ::= (l, t, \iota_1 \dots \iota_n, \chi)$
	(function) $F ::= (x_1, \dots, x_n)B_1, \dots, B_m$
	(program) $P ::= F_1, \dots, F_n$

Fig. 2. CIR Syntax

ers can. The template names on the blocks within a container are used to group blocks into templates. Specifically, a template t is the list of blocks with template name t .¹ Executable functions use containers to generate code dynamically by *copying templates* (t) into *code regions* (r) and *filling holes* ($[h]$). We call a copy of a template an *instance* (p). Holes ($[h]$) allow generating code to specialize instances. Generating code can dictate the control flow of generated code by filling destination holes, and can insert values into the generated instructions by filling value holes ($x := [h]$).

The instruction $r := \underline{\text{start}}\ l$ allocates a code region and binds it to r . All templates copied into r should come from the (container) function named l . The instruction $p := \underline{\text{copy}}\ t\ \underline{\text{into}}\ r$ puts an instance of t after the instances previously copied into r .² p becomes a reference to the new instance. The two fill instructions ($\underline{\text{fill}}\ p.[h]\ \underline{\text{with}}\ \dots$) both fill the hole h of the specified instance (p). The second form allows one instance to refer directly to a label in another instance. That way, the generating code can put inter-template jumps in the generated code. Finally, $x := \underline{\text{end}}\ r$ completes the generation of a function and binds x to the resulting function pointer.

Inter-template jumps merit further discussion. It makes no sense for a template to refer directly to a label in another template. To which instance of a template would such a label refer? Therefore, the generating code must create all inter-template control flow using the $\underline{\text{fill}}\ p_1.[h]\ \underline{\text{with}}\ p_2.l$ instruction. There is an important exception: If a template's last block has a fall-through destination (such as $\underline{\text{jmp}}\ \circ$), then control will flow from an instance of that template to the next instance in the code region. That is, fall-through destinations refer to the code region's next block, not the container function's next block.

¹ The order of blocks in the template is the same as their order in the function.

Template names within executable functions are present only to simplify the syntax.

² The implementation detects buffer overflow and moves the region to a larger buffer.

<pre> (u) (<i>dot_gen</i>, \rightarrow, <i>r</i> := <u>start</u> <i>dot</i> <i>p</i>₁ := <u>copy</u> <i>t</i>₁ <u>into</u> <i>r</i> <i>i</i> := 0, <u>jmp</u> \circ) (<i>l</i>₃, \rightarrow, <i>tst</i> := <i>i</i> < <i>size</i>(<i>u</i>), <u>jnz</u> <i>tst</i>?<i>l</i>₅ : \circ) (<i>l</i>₄, \rightarrow, <i>p</i>₂ := <u>copy</u> <i>t</i>₂ <u>into</u> <i>r</i> <u>fill</u> <i>p</i>₂.<i>h</i>₁ <u>with</u> <i>u</i>[<i>i</i>] <u>fill</u> <i>p</i>₂.<i>h</i>₂ <u>with</u> <i>i</i> <i>i</i> := <i>i</i> + 1, <u>jmp</u> <i>l</i>₃) (<i>l</i>₅, \rightarrow, <i>p</i>₃ := <u>copy</u> <i>t</i>₃ <u>into</u> <i>r</i> <i>f</i> := <u>end</u> <i>r</i>, <u>retn</u> <i>f</i>) </pre>	<pre> (v) (<i>dot</i>, <i>t</i>₁, <i>result</i> := 0, <u>jmp</u> \circ) (<i>l</i>₁, <i>t</i>₂, <i>x</i> := [<i>h</i>₁] <i>i</i> := [<i>h</i>₂] <i>tmp</i> := <i>x</i> * <i>v</i>[<i>i</i>] <i>result</i> := <i>result</i> + <i>tmp</i>, <u>jmp</u> \circ) (<i>l</i>₂, <i>t</i>₃, , <u>retn</u> <i>result</i>) </pre>
--	---

Fig. 3. Translation of the `dot_gen` example (simplified)

Figure 3 shows our `dot_gen` example in CIR (extended with arrays, etc.). The translation of Section 5 would introduce more variables but would otherwise produce similar output. Notice the function `dot` contains holes, so it is not executable. When `dot_gen` is called, it creates a code region `r` in which to generate a function. It then copies `dot`'s prologue, `t1`, and enters the loop consisting of blocks `l3` and `l4`. The loop copies an instance of `t2` and binds `p2` to it. `p2` is used when filling holes `h1` and `h2` to specify which instance's holes to fill. Hole `h1` is filled with the value in `u[i]` and hole `h2` is filled with the counter `i`. After the loop, `dot_gen` copies `dot`'s epilogue, `t3`, and ends the code generation.

To employ inter-template optimizations on `dot`, it helps to have an approximation of the order its templates might be copied. Put another way, we seek to approximate what the fall-through destinations of `t1` and `t2` might refer to. Without looking at `dot_gen`, we must pessimistically assume that instances of `t1` and `t2` might be followed by instances of `t1`, `t2`, or `t3`, but an analysis of `dot_gen` reveals that an instance of `t1` would never follow an instance of `t1` or `t2`.

5 Translation

This section describes a translation from Mini-Cyclone (Section 3) to CIR (Section 4). We begin with some highlights.

Translating a lexically scoped language with RTCG into an unstructured block-based language is challenging because the source language intertwines generated and generating code. In the target language, a function is either executable or a template container, never a mixture. The one-pass top-down translation we present untangles the levels by maintaining an environment for each level.

If a function `f` contains “`codegen g . . .`,” then we call `f` the *parent* of `g` and `g` a *child* of `f`. During the translation, exactly one function is *active*; instructions

are emitted into the active function. Entering a `cut` activates the active function’s parent, and entering a `splice` activates the active function’s child. Entering a `codegen` is like entering a `splice` except that it spawns a new child and then activates the child. For `fill e`, we emit a value hole and then activate the parent just long enough to translate e and emit a `fill` instruction. Notice that an explicit Mini-Cyclone `fill` always corresponds to a value hole.

All the code in a template must be copied into a code region at once. Consequently, whenever two pieces of code might not be copied consecutively, we must introduce a new template. Our translation conservatively creates new templates on each transition from parent to child (when entering a `splice` or `codegen`, and when leaving a `cut`).

Multiple templates create the possibility for inter-template jumps. To see how such control flow arises, consider code of the form, “`while(e){ cut{ ... } } ...`” The `cut` causes the code following the loop and the loop guard to be in different templates. Hence, the control transfer for a false guard is inter-template. If the loop body had not contained a `cut`, then no inter-template jump would have been necessary. After translating the loop body, the translation checks if the template following the loop is the same as the template for the guard. If not, then the translation places a destination hole in the guard, and issues an inter-template fill (`fill p1.[h] with p2.l`).

The above difficulties make translating RTCG into an unstructured block-based language inherently more complex than standard translations. To explore all the subtleties, we present a full translation from Mini-Cyclone into CIR. In order to emphasize the interesting aspects, some parts are informal; a formal presentation of the translation as a pure monadic-style function over a compilation environment is available [15].

5.1 Translation Environment and Preliminaries

The translation of a Mini-Cyclone term is with respect to a global environment, Φ , that contains local environments, \mathcal{E} , for each function being simultaneously translated. A global environment is a (functional) record of type τ_Φ where a local environment has type $\tau_\mathcal{E}$.

$$\tau_\Phi = \{\mathbf{parents} : \tau_\mathcal{E} \text{ list}, \mathbf{active} : \tau_\mathcal{E}, \mathbf{children} : \tau_\mathcal{E} \text{ list}\}$$

A global environment represents a sequence of the form:

$$\underbrace{\mathcal{E}_1, \dots, \mathcal{E}_{k-1}}_{\mathbf{parents}}, \underbrace{\mathcal{E}_k}_{\mathbf{active}}, \underbrace{\mathcal{E}_{k+1}, \dots, \mathcal{E}_n}_{\mathbf{children}}$$

The head of the list **parents** is \mathcal{E}_{k-1} and the head of the list **children** is \mathcal{E}_{k+1} . In general, if f is the parent of g , then f immediately precedes g in the sequence. We need a list because `codegen` expressions can nest.

We treat $\tau_\mathcal{E}$ informally because it has the same information as an environment for a conventional compiler: It is a partially generated target function into

which we (imperatively) add blocks. It also contains a distinguished *current block* into which we (imperatively) emit instructions. When we are done generating a function, we use *addFun* \mathcal{E} to add it to the output.

A tricky aspect of any translation is relating source variables to target variables while easily generating new names. We finesse this issue by asserting:

$$\text{VAR} = \text{CYCVAR} \cup \text{GENVAR} \quad \text{CYCVAR} \cap \text{GENVAR} = \emptyset$$

Hence variables in GENVAR cannot clash with names in the source program. We use the function *newVar* to generate fresh elements of GENVAR.

We assume that the base syntactic classes are isomorphic and that we have functions witnessing these isomorphisms: Given an element of any class (call it a), we can get the corresponding variable ($\mathcal{V}(a)$), label ($\mathcal{L}(a)$), hole ($\mathcal{H}(a)$), template name ($\mathcal{T}(a)$), code region ($\mathcal{R}(a)$), or instance ($\mathcal{P}(a)$). This technical trick helps us exploit subtle invariants. For example, the translation has the property that only one instance of a template is live at a time, so we use the template name (t) to induce an instance name ($\mathcal{P}(t)$).

5.2 Non-RTCG Translation

Each top-level function of the input program is translated independently, so to compile a program, we simply compile each source function and collect all the target functions created.

The translation of a function, $[[D]]$, takes an environment, Φ , as an argument. For top-level functions, Φ should be the empty sequence of function environments. For a D declared in **codegen**, Φ .**parents** will begin with D 's parent (see the next section). The role of $[[D]]$ is to translate the function's body under an environment where Φ .**active** is a new function environment. Roughly, *newFun* creates a function environment for $f(x_1, \dots, x_n) s$ in which there is one empty block with label $\mathcal{L}(f)$ and template name $\mathcal{T}(f)$. Using an informal notation:

$$\begin{aligned} [[f(x_1, \dots, x_n) s]] \Phi = & \\ & \mathbf{let} \ \mathcal{E} = \mathit{newFun} \ (f(x_1, \dots, x_n) \ s) \ \mathbf{in} \\ & [[s]] \ \Phi[\mathbf{active} = \mathcal{E}]; \\ & \mathit{addFun} \ \mathcal{E} \end{aligned}$$

The main effect of $[[D]]$ is to add a target function that is the translation of D . (As a side-effect, any **codegen** terms in the body of D give rise to target functions too.) Notice we use the notation $\Phi[\mathbf{label} = F]$ for functional-record update; the result is a record with the same fields as Φ except field **label** has value F .

The translation of statements, $[[s]]$, also takes a Φ . The translation of expressions, $[[e]]$, leaves its result in a CIR variable provided as input, so it takes a Φ and an x . The definitions $[[s]]$ and $[[e]]$ are inductive over the structure of Mini-Cyclone statements and expressions. The rest of this section presents the non-RTCG cases for these definitions, beginning with the simpler cases.

All of the non-RTCG expression forms appear straightforward to translate. Most of the work lies in generating fresh variables to hold intermediate results. The term $emit \ \Phi \ \iota$ updates Φ such that ι is appended to the end of the active function's current block. Analogously, $emit \ \Phi \ \chi$ replaces the control transfer of the current block with χ . Finally, $e_1; e_2$ means, "Do e_1 then e_2 ."

$$\begin{array}{lll}
[[i]] \ \Phi \ x = & [[e_1 + e_2]] \ \Phi \ x = & [[e_0(e_1, \dots, e_n)]] \ \Phi \ x = \\
\quad emit \ \Phi \ x := i & \quad \mathbf{let} \ x_1 = newVar \ \mathbf{in} & \quad \mathbf{let} \ x_0 = newVar \ \mathbf{in} \\
& \quad \mathbf{let} \ x_2 = newVar \ \mathbf{in} & \quad \dots \\
[[f]] \ \Phi \ x = & [[e_1]] \ \Phi \ x_1; & \quad \mathbf{let} \ x_n = newVar \ \mathbf{in} \\
\quad emit \ \Phi \ x := \mathcal{L}(f) & [[e_2]] \ \Phi \ x_2; & \quad [[e_0]] \ \Phi \ x_0; \\
& \quad emit \ \Phi \ x := x_1 + x_2 & \quad \dots \\
[[x]] \ \Phi \ y = & & [[e_n]] \ \Phi \ x_n; \\
\quad emit \ \Phi \ y := x & & \quad emit \ \Phi \ x := x_0(x_1, \dots, x_n)
\end{array}$$

The following statement cases are just as straightforward:

$$\begin{array}{lll}
[[x := e]] \ \Phi = & [[s_1; s_2]] \ \Phi = & [[\mathbf{return} \ e]] \ \Phi = \\
\quad [[e]] \ \Phi \ x & \quad [[s_1]] \ \Phi; & \quad \mathbf{let} \ x = newVar \ \mathbf{in} \\
& \quad [[s_2]] \ \Phi & \quad [[e]] \ \Phi \ x; \\
& & \quad emit \ \Phi \ \underline{\mathbf{retn}} \ x
\end{array}$$

The control-flow constructs (**if** and **while**) are more complicated. Both definitions follow the same strategy: Generate the subterms while remembering the entry and exit blocks for each. Then replace various transfers to create the correct control flow. This order is important for RTCG. As explained earlier, the translation may have to create inter-template control-flow at these points. We use the term $genDest$ for this purpose. In the absence of RTCG, $genDest(b_{src}, b_{dst}) = b_{dst}$; the next section gives a complete definition. The term $changeBlock$ puts a new block in the active function and makes the new block the current one. $changeBlock$ returns a pair of the old current block's label and the new current block's label. The new block has the same template name as the old block. The term $emit \ \Phi \ b: \chi$ replaces the transfer in block b with χ .

$$\begin{array}{ll}
[[\mathbf{if}(e) \ s_1 \ \mathbf{else} \ s_2]] \ \Phi = & [[\mathbf{while}(e) \ s]] \ \Phi = \\
\quad \mathbf{let} \ x = newVar \ \mathbf{in} & \quad \mathbf{let} \ (b_0, b_t) = changeBlock \ \Phi \ \mathbf{in} \\
\quad [[e]] \ \Phi \ x; & \quad \mathbf{let} \ x = newVar \ \mathbf{in} \\
\quad \mathbf{let} \ (b_0, b_t) = changeBlock \ \Phi \ \mathbf{in} & \quad [[e]] \ \Phi \ x; \\
\quad [[s_1]] \ \Phi; & \quad \mathbf{let} \ (b_1, b_b) = changeBlock \ \Phi \ \mathbf{in} \\
\quad \mathbf{let} \ (b_1, b_f) = changeBlock \ \Phi \ \mathbf{in} & \quad [[s]] \ \Phi; \\
\quad [[s_2]] \ \Phi; & \quad \mathbf{let} \ (b_2, b_e) = changeBlock \ \Phi \ \mathbf{in} \\
\quad \mathbf{let} \ (b_2, b_m) = changeBlock \ \Phi \ \mathbf{in} & \quad \mathbf{let} \ d_e = genDest \ \Phi \ (b_1, b_e) \ \mathbf{in} \\
\quad \mathbf{let} \ d_f = genDest \ \Phi \ (b_0, b_f) \ \mathbf{in} & \quad \mathbf{let} \ d_t = genDest \ \Phi \ (b_2, b_t) \ \mathbf{in} \\
\quad \mathbf{let} \ d_m = genDest \ \Phi \ (b_1, b_m) \ \mathbf{in} & \quad emit \ \Phi \ b_0: \underline{\mathbf{jmp}} \ \circ; \\
\quad emit \ \Phi \ b_0: \underline{\mathbf{jnz}} \ x ? \circ : d_f; & \quad emit \ \Phi \ b_1: \underline{\mathbf{jnz}} \ x ? \circ : d_e; \\
\quad emit \ \Phi \ b_1: \underline{\mathbf{jmp}} \ d_m; & \quad emit \ \Phi \ b_2: \underline{\mathbf{jmp}} \ d_t \\
\quad emit \ \Phi \ b_2: \underline{\mathbf{jmp}} \ \circ
\end{array}$$

5.3 RTCG Translation

We now explain how the four Mini-Cyclone constructs specific to RTCG are compiled. We already have all the machinery we need to translate `codegen`:

$$\begin{aligned}
[[\text{codegen } f(x_1, \dots, x_n)s]] \Phi \ x = & \\
\text{emit } \Phi \ \mathcal{R}(f) := \underline{\text{start}} \ \mathcal{L}(f); & \\
\text{emit } \Phi \ \mathcal{P}(f) := \underline{\text{copy}} \ \mathcal{T}(f) \ \underline{\text{into}} \ \mathcal{R}(f); & \\
[[f(x_1, \dots, x_n)s]] \Phi [\text{parents} = \Phi.\text{active} :: \Phi.\text{parents}]; & \\
\text{emit } \Phi \ x := \underline{\text{end}} \ \mathcal{R}(f) &
\end{aligned}$$

The translation creates both generating code and generated code. At the generating level, we first allocate the code region ($\mathcal{R}(f) := \underline{\text{start}} \ \mathcal{L}(f)$) and copy the child's first template ($\mathcal{P}(f) := \underline{\text{copy}} \ \mathcal{T}(f) \ \underline{\text{into}} \ \mathcal{R}(f)$). For translating the child, we provide an environment where the current active function is the parent. That way, `cut` statements in s will use the correct local environment. Recall that $[[D]]$ takes care of creating a new function (the child) and considering it active. The translation of s will add code at various levels. After the child has been translated, the parent (that is, the active function of Φ) will not have any more instructions emitted into it that manipulate this child. So the last step is to emit code that converts the code region to executable code ($x := \underline{\text{end}} \ \mathcal{R}(f)$). Notice that we use the original Φ except when translating the child function.

For translating the remaining constructs, special notation for shifting the active function will prove useful:

$$\begin{aligned}
\uparrow \Phi &= \{\text{parents} = \Phi.\text{active} :: \Phi.\text{parents}, \ \text{active} = \text{head}(\Phi.\text{children}), \\
&\quad \text{children} = \text{tail}(\Phi.\text{children})\} \\
\downarrow \Phi &= \{\text{children} = \Phi.\text{active} :: \Phi.\text{children}, \ \text{active} = \text{head}(\Phi.\text{parents}), \\
&\quad \text{parents} = \text{tail}(\Phi.\text{parents})\}
\end{aligned}$$

So $\uparrow \Phi$ makes the child active and $\downarrow \Phi$ makes the parent active. Translation of well-formed source code will never apply *tail* to an empty list.

The translations for `cut` and `splice` are pleasingly symmetric:

$$\begin{aligned}
[[\text{cut } s]] \Phi = & \quad [[\text{splice } s]] \Phi = \\
\text{emit } \Phi \ \underline{\text{jmp}} \ \circ; & \quad \text{newTemplate } (\uparrow \Phi); \\
[[s]] (\downarrow \Phi); & \quad [[s]] (\uparrow \Phi); \\
\text{newTemplate } \Phi & \quad \text{emit } (\uparrow \Phi) \ \underline{\text{jmp}} \ \circ
\end{aligned}$$

Recall that the purpose of `cut` s is to execute s in the parent of the function in which the `cut` appears. Dually, `splice` s puts the code fragment s in the child of the function in which the `splice` appears. In both cases, we have a child-to-parent transition (beginning of the cut, end of the splice) and a parent-to-child transition (end of the cut, beginning of the splice). For child-to-parent transitions, we end the current block so that control will flow to a new template that an ensuing parent-to-child transition creates. For parent-to-child transitions, we put a new template in the child because the parent may choose to copy the code after the transition at different times than the code before the transition. The

auxiliary term *newTemplate* creates a new template in the child and a copy instruction in the parent. (*activeFunction* Φ retrieves the active function's name. *setTemplateOfCurrent* Φ t changes the current block's template name to t .)

```

newTemplate  $\Phi$  =
  let  $f$  = activeFunction  $\Phi$  in
  let  $x$  = newVar in
  let  $(b_0, b_1)$  = changeBlock  $\Phi$  in
  setTemplateOfCurrent  $\Phi$   $T(x)$ ;
  emit ( $\downarrow$   $\Phi$ )  $\mathcal{P}(x) := \underline{\text{copy}}\ T(x)\ \underline{\text{into}}\ \mathcal{R}(f)$ 

```

The translation of **fill** e is straightforward at this point. In the parent, we translate e and emit a **fill** instruction for a new hole. In the child, we emit a value hole. (*currentTemplate* retrieves the current block's template name.)

```

[[fill  $e$ ]]  $\Phi$   $x$  =
  let  $x_1$  = newVar in
  [[ $e$ ]] ( $\downarrow$   $\Phi$ )  $x_1$ ;
  let  $x_2$  = newVar in
  let  $t$  = currentTemplate  $\Phi$  in
  emit ( $\downarrow$   $\Phi$ ) fill  $\mathcal{P}(t).$ [ $\mathcal{H}(x_2)$ ] with  $x_1$ ;
  emit  $\Phi$   $x :=$  [ $\mathcal{H}(x_2)$ ]

```

So far, the translation does not seem to use any **fill** $p_1.[h]$ **with** $p_2.l$ instructions. As discussed in Section 4, we need such instructions for inter-template jumps. We have just seen that new templates are created when **cut** statements occur inside of other source constructs. This fact is precisely why the control-flow jumps in the translation of **if** and **while** may be inter-template. We now have the machinery to define easily *genDest* for the general case: We simply check whether the templates of the source and destination blocks are the same (using *templateOf* Φ b to retrieve the template name in the block with label b). If the same, we just return the destination's label. If different, we create a hole, emit a fill for the hole in the parent, and return the hole.

```

genDest  $\Phi$  ( $b_{src}, b_{dst}$ ) =
  let  $t_{src} =$  templateOf  $\Phi$   $b_{src}$  in
  let  $t_{dst} =$  templateOf  $\Phi$   $b_{dst}$  in
  if  $t_{src} = t_{dst}$  then  $b_{dst}$ 
  else (let  $x =$  newVar in
        emit ( $\downarrow$   $\Phi$ ) fill  $\mathcal{P}(t_{src}).$ [ $\mathcal{H}(x)$ ] with  $\mathcal{P}(t_{dst}).b_{dst}$ ;
        [ $\mathcal{H}(x)$ ])

```

6 Optimization

A goal of our research is to study the interaction between RTCG and optimization. By careful design of CIR, we have made the application of traditional data-flow optimizations straightforward. These optimizations can be split into an analysis and a transformation phase. We consider these phases individually.

6.1 Analysis

A preliminary step in data-flow analyses is to build a control-flow graph (CFG) that conservatively approximates the dynamic control flow of the function. For a conventional language, we just put an edge between a block and the blocks that its control transfer mentions. In the presence of RTCG with holes or fall-throughs for destinations, control transfers are no longer apparent. In this section, we describe how to compute an appropriate CFG for any CIR function.

A CFG for a function F , which we write $\text{CFG}(F)$, is a directed graph where the nodes are the labels in F . To compute $\text{CFG}(F)$, we start with an empty graph and add edges. Adding edges for intra-template control transfers is conventional. For inter-template control transfers, we need to approximate how generating code might create control flow among F 's templates. Such control flow can result from filling holes in transfer instructions. It can also result from copying templates that end with fall-through destinations. One conservative approach would be to assume that a hole destination could be filled with any label of any instance of F and that F 's templates could be copied in any order.

By analyzing the parent of F (that is, the function containing $r := \text{start } f$ where f is the name of F), we can compute a less conservative CFG. Our translation has the property that executable functions have no parents and non-executable functions have exactly one parent. We denote F 's parent by $F \downarrow$. We explain how to extend the analysis to multiple parents at the end of the section.

We begin with the intra-template edges. Let $F = (x_1, \dots, x_n)B_1, \dots, B_m$ and $B_i = (l_i, t_i, \iota_{i1} \dots \iota_{ip}, \chi_i)$. The intra-template destinations for block B_i are:

$$\begin{aligned} \text{intra}_i([h]) &= \emptyset \\ \text{intra}_i(l_j) &= \begin{cases} \{l_j\} & \text{if } t_i = t_j \\ \text{undefined} & \text{otherwise} \end{cases} \\ \text{intra}_i(\circ) &= \begin{cases} \{l_{i+1}\} & \text{if } i < m \text{ and } t_i = t_{i+1} \\ \emptyset & \text{otherwise} \end{cases} \end{aligned}$$

Hole destinations are handled by the inter-template part of the algorithm (even though a parent may fill the hole with a label in the same instance). An explicit destination must be to a block in the same template, a fact the undefined case emphasizes. A fall-through destination is an intra-template transfer if and only if the next block in the function belongs to the same template. We extend the definition to transfers ($\text{intra}_i(\chi)$) and functions ($\text{intra}(F)$) by taking the union of $\text{intra}_i(d)$ over all constituent destinations, d .

For functions without parents (executable functions), $\text{CFG}(F) = \text{intra}(F)$. For those with parents (container functions), $\text{CFG}(F) = \text{intra}(F) \cup \text{inter}(F)$. So we now turn to $\text{inter}(F)$, assuming we have $\text{CFG}(F \downarrow)$. That is, we compute the parent's CFG and use it to compute the child's CFG. For now, assume that $F \downarrow$ contains exactly one $r := \text{start } l$ instruction. Our goal is to determine what edges we must add to $\text{CFG}(F)$ because of the way $F \downarrow$ manipulates the templates. We assume every instruction in $F \downarrow$ is reachable, so the edges added

are the union of the edges each instruction adds:

$$\mathbf{inter}(F) = \bigcup_{\iota \in F \downarrow} \mathbf{inter}(\iota)$$

Instructions add edges by filling holes and copying templates:

$$\mathbf{inter}(\iota) = \begin{cases} \{(l_i, l_j)\} & \text{if } \iota = \mathbf{fill} \ p_1.[h] \ \mathbf{with} \ p_2.l_j \text{ and } h \text{ is in } \chi_i \\ \mathbf{copypred}(\iota) \times \{l\} & \text{if } \iota \text{ copies } t \text{ and } t\text{'s first block has label } l \\ \emptyset & \text{otherwise} \end{cases}$$

The first case includes edges introduced by filling holes. The second case is for fall-through destinations: If ι copies t , then consider all t' that could have been the most recently copied template when control reaches ι . If the last block of t' has a fall-through destination, then we must add an edge from the last block of t' to the first block of t . We write $\mathbf{copypred}(\iota)$ for the set of the labels of such “last blocks”. Determining $\mathbf{copypred}(\iota)$ for each ι in $F \downarrow$ is a data-flow problem. We define the function $\mathbf{fallthru}(t)$ for a template in function F to be the label of the last block if that block has a fall-through destination:

$$\mathbf{fallthru}(t) = \begin{cases} \{l\} & \text{if } (l, t, \dots, \chi) \text{ is } t\text{'s last block and } \chi \text{ contains } \circ \\ \emptyset & \text{otherwise} \end{cases}$$

For an instruction sequence, \mathbf{last} denotes the last template copied, if any:

$$\mathbf{last}(\iota_1; \dots; \iota_n) = \begin{cases} \{t\} & \text{if } \iota_n = p := \mathbf{copy} \ t \ \mathbf{into} \ r \\ \mathbf{last}(\iota_1; \dots; \iota_{n-1}) & \text{otherwise} \end{cases}$$

$$\mathbf{last}(\cdot) = \emptyset$$

We extend \mathbf{last} to blocks via $\mathbf{last}(l, t, \iota_1, \dots, \iota_n, \chi) = \mathbf{last}(\iota_1, \dots, \iota_n)$. We can now state the relevant data-flow equations.

$$\mathbf{in}(B) = \bigcup_{B' \in \mathbf{pred}(B)} \mathbf{out}(B') \quad \mathbf{out}(B) = \begin{cases} \mathbf{fallthru}(t) & \text{if } \mathbf{last}(B) = \{t\} \\ \mathbf{in}(B) & \text{if } \mathbf{last}(B) = \emptyset \end{cases}$$

These equations are solvable using standard techniques [1, 12] with $\mathbf{in}(B) = \emptyset$ for the first block of $F \downarrow$. Note that the definition of $\mathbf{in}(B)$ is where we assume we have already determined $\mathbf{CFG}(F \downarrow)$. Using the data-flow solution and letting $\iota_{(l,i)}$ denote the i^{th} instruction in block l , we can define $\mathbf{copypred}$:

$$\mathbf{copypred}(\iota_{(l,i)}) = \begin{cases} \mathbf{in}(B) & \text{if } i = 1 \text{ and } B = (l, \dots) \\ \mathbf{fallthru}(t) & \text{if } \iota_{(l,i-1)} = p := \mathbf{copy} \ t \ \mathbf{into} \ r \\ \mathbf{copypred}(l, i-1) & \text{otherwise} \end{cases}$$

We now discuss how to generalize the CFG analysis. If F has multiple parents, we just add the edges that each parent requires. Similarly, if one parent uses multiple code regions, we just repeat the process for each one — because code

regions and instance pointers are second-class constructs, we always know what code region an instruction affects. The interesting extension is for when functions are mutually generating. That is, nothing in CIR prevents a cycle in the graph in which the edge (F, G) means F is a parent for G . Here we cannot compute the CFG for a parent before a child. The solution is to iterate the program-wide analysis in a manner just like conventional inter-procedural analysis. Note that our translation never produces mutually generating functions.

6.2 Transformation

Although analyses are little affected by RTCG, transformations require care. For example, dead-code elimination might eliminate code containing a hole. Doing so requires removing the attendant fills in the parent. Similarly, if a hole is duplicated, the fills must be duplicated. Our compiler summarizes the relevant transformation effects and then uses this summary to update the parent.

Unfortunately, treating holes as constants can be too aggressive: Each time a template is copied, its holes may be filled with *different* values. For instance, the block l_1 is in a loop in $\text{CFG}(\text{dot})$ (Figure 3) and i (if holes are constants) is invariant. However, hoisting the instruction $i := [h_2]$ out of the loop would not preserve the program’s semantics. The hoist would have been valid had the source and destination of the transfer been in the same instance (equivalently, if the edge is in $\text{intra}(\text{dot})$). A principled understanding of transformations will require having two kinds of edges in CFG – those from **intra** and those from **inter**. For now, our compiler does not move holes across templates.

Note that a sound transformation of a parent will not require updating a child: Because the parent treats the child as data, any correct optimization will copy the templates in the same order. However, optimizing a parent may reveal a less conservative CFG for a child, thus enabling child optimizations.

7 Related Work

Compared to other RTCG systems, Cyclone has a unique combination of features: type-safety, a simple operational model, and user-level control. 'C [5, 14] is an unsafe language providing rich low-level control over RTCG. Users can choose among different run-time compilers to trade code quality for compilation time. In essence, more costly compilers delay inter-template optimizations until run time. MetaML [17, 16] and PML [18, 4] provide type safety and user-level control but do not have an operational model matching the underlying machine. Therefore, the costs and benefits of RTCG are harder to predict. Fabius [10, 9], Tempo [13, 3], and DyC [2, 7, 6] all apply RTCG automatically based on user annotations, thus relieving the user of low-level control.

Our contribution is a formal description of compilation and optimization for template-based RTCG. Prior work has used ad hoc techniques. Tempo produces template C code that is compiled to object code by GCC. From the resulting object code, the Tempo run-time system extracts templates. This approach is

easier to implement, but is less robust, than writing a full compiler. DyC is an impressive optimizing compiler with pervasive support for RTCG. Most data-flow optimizations are performed *before* the program is staged, so it is necessary only to preserve the user annotations during optimization. This preservation is sometimes difficult, so some optimizations are disabled [6].

8 Conclusion

In this paper, we have presented a formal translation from a language (Mini-Cyclone) with explicit RTCG to an intermediate language (CIR), and we have shown how to apply standard data-flow optimizations. Hopefully, formalizing these techniques will allow other compilers to integrate RTCG more easily.

We have implemented an optimizing compiler for a realistic source language with RTCG; this paper distills the language and compiler down to its essence. Our system currently performs only a small set of optimizations, notably register allocation and null-check elimination. Nonetheless, initial results are promising: For an implementation of RSA, we applied RTCG to code accounting for 30% of the execution time and obtained a 10% speedup over the unspecialized program.

References

1. Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. *Compilers, Principles, Techniques and Tools*. Addison-Wesley, 1986.
2. Joel Auslander, Matthai Philipose, Craig Chambers, Susan J. Eggers, and Brian N. Bershad. Fast, effective dynamic compilation. In *ACM Conference on Programming Language Design and Implementation*, pages 149–159, May 1996.
3. Charles Consel and François Noël. A general approach to run-time specialization and its application to C. In *Twenty-Third ACM Symposium on Principles of Programming Languages*, pages 145–156, January 1996.
4. Rowan Davies and Frank Pfenning. A modal analysis of staged computation. In *Twenty-Third ACM Symposium on Principles of Programming Languages*, pages 258–270, January 1996.
5. Dawson Engler, Wilson Hsieh, and M. Frans Kaashoek. 'C: A language for fast, efficient, high-level dynamic code generation. In *Twenty-Third ACM Symposium on Principles of Programming Languages*, pages 131–144, January 1996.
6. Brian Grant, Markus Mock, Matthai Philipose, Craig Chambers, and Susan J. Eggers. DyC: An expressive annotation-directed dynamic compiler for C. Technical Report UW-CSE-97-03-03, Department of Computer Science and Engineering, University of Washington. Updated May 12, 1999.
7. Brian Grant, Matthai Philipose, Markus Mock, Craig Chambers, and Susan J. Eggers. An evaluation of staged run-time optimizations in DyC. In *ACM Conference on Programming Language Design and Implementation*, pages 293–304, May 1999.
8. Luke Hornof and Trevor Jim. Certifying compilation and run-time code generation. In *ACM Conference on Partial Evaluation and Semantics-Based Program Manipulation*, pages 60–74, January 1999.
9. Peter Lee and Mark Leone. Optimizing ML with run-time code generation. In *ACM Conference on Programming Language Design and Implementation*, pages 137–148, May 1996.

10. Mark Leone and Peter Lee. Lightweight run-time code generation. In *ACM Conference on Partial Evaluation and Semantics-Based Program Manipulation*, pages 97–106, June 1994.
11. Greg Morrisett, David Walker, Karl Cray, and Neal Glew. From System F to typed assembly language. *ACM Transactions on Programming Languages and Systems*, 21(3):528–569, May 1999.
12. Steven S. Muchnick. *Advanced Compiler Design and Implementation*. Morgan Kaufmann Publishers, 1997.
13. François Noël, Luke Hornof, Charles Consel, and Julia L. Lawall. Automatic, template-based run-time specialization: Implementation and experimental study. Technical Report 1065, Institut de Recherche en Informatique et Systèmes Aléatoires (IRISA), November 1996.
14. Massimiliano Poletto, Dawson Engler, and M. Frans Kaashoek. tcc: A system for fast, flexible and high-level dynamic code generation. In *ACM Conference on Programming Language Design and Implementation*, pages 109–121, June 1997.
15. Frederick Smith, Dan Grossman, Greg Morrisett, Luke Hornof, and Trevor Jim. Compiling for runtime code generation (extended version). Technical report, Department of Computer Science, Cornell University, October 2000.
16. Walid Taha. A sound reduction semantics for untyped CBN mutli-stage computation. Or, the theory of MetaML is non-trivial (extended abstract). In *ACM Conference on Partial Evaluation and Semantics-Based Program Manipulation*, pages 34–43, January 2000.
17. Walid Taha and Tim Sheard. Multi-stage programming with explicit annotations. In *ACM Conference on Partial Evaluation and Semantics-Based Program Manipulation*, pages 203–217, June 1997.
18. Philip Wickline, Peter Lee, and Frank Pfenning. Run-time code generation and Modal-ML. In *ACM Conference on Programming Language Design and Implementation*, pages 224–235, June 1998.