

# An Extensible, Higher-Order Intermediate Representation

ANONYMOUS AUTHOR(S)

Traditional compilers, designed for optimizing low-level code, fall short when dealing with modern, computation-heavy applications like image processing, machine learning, or numerical simulations. Optimizations should understand the primitive operations of the specific application domain and thus happen on that level.

Domain-specific languages (DSLs) fulfill these requirements. However, DSL compilers reinvent the wheel over and over again as standard optimizations, code generators, and general infrastructure & boilerplate code must be reimplemented for each DSL compiler.

This paper presents THORIN, an extensible, higher-order intermediate representation. At its core, THORIN is a pure type system and, hence, a form of a typed lambda calculus. Developers can declare the signatures of new (domain-specific) operations, called *axioms*. An axiom can be the declaration of a function, a type operator, or any other entity with a possibly polymorphic, polytypic, and/or dependent type. This way, developers can extend THORIN at any low or high level and bundle them in a *plugin*. Plugins extend the compiler and take care of optimizing and lowering the plugin’s axioms.

We show the expressiveness and effectiveness of THORIN in three case studies: Low-level plugins that operate at the same level of abstraction as LLVM, a regular-expression matching plugin, and an automatic differentiation pass. We show that in all three studies, THORIN produces code that has state-of-the-art performance.

## 1 INTRODUCTION

Dennard scaling enabled the continuous growth of single-thread performance of legacy code for decades. After its decline in the early 2000s, domain-specific languages (DSLs) have received a lot of attention in the effort to harvest as much performance as possible in a productive way. DSLs offer specialized abstractions tailored to specific problem domains enhancing the programmer’s productivity on one side and enabling the generation of high-performance code on the other side. To generate actual code, DSL compilers typically resort to existing compiler frameworks such as LLVM [30]. However, there is a significant gap in the abstractions DSLs provide and the low-level nature of these compiler frameworks which often leads to the situation that DSL compilers have an additional, more high-level intermediate representation (IR) that bridges the gap between the DSL and the back-end compiler.

This more high-level IR is often designed and implemented for each new DSL from scratch. Current research in high-level IRs, most notably MLIR [31], seeks to provide a least common denominator for the basis of such an IR to facilitate the reuse of core data structures and algorithms across different DSLs but not much more. While MLIR provides extensibility to host domain-specific *dialects*, it is still low-level in several aspects. First, it relies on first-order control flow based on control-flow graphs (CFGs), thereby relying on multiple concepts to represent control flow including functions, basic blocks, instructions, and so-called regions to delineate code for high-level transformations. Second, MLIR does not provide a type system that is expressive enough to host the type systems of DSLs. It relies on the dialects to implement their own type systems in C++. This causes more implementation effort for the DSL implementer and an unclear notion of well-typedness when multiple dialects interact. Finally, because MLIR is designed with “little builtin and everything customizable” [31], it is very hard to give a formal account of its static and dynamic semantics that would encompass all DSLs. The type system therefore can give fewer guarantees regarding the semantics of and between individual DSLs.

In this paper, we want to go one step further and present THORIN, a higher-order IR that provides a type system that sets out to be expressive enough to host a wide range of DSLs. Consider the

following DSL interface for a *matrix* plugin in THORIN’s *surface language*—the textual interface of THORIN that includes enough syntactic sugar to specify complex types:

```
.ax %matrix.Mat:  $\Pi$  [n m: .Nat, T: *]  $\rightarrow$  *; // n  $\times$  m matrix with element type T
.ax %matrix.zip:  $\Pi$ . [n m: .Nat, T: *] [f: [T, T]  $\rightarrow$  T] [a b: %matrix.Mat (n, m, T)]
     $\rightarrow$  %matrix.Mat (n, m, T);
```

The so-called *axiom* `%matrix.Mat` is a type constructor that expects two dimensions `n` and `m` and an element type `T` to construct a matrix. The `%matrix.zip` axiom is polymorphic in the arguments’ dimensions and element type (whose values will be inferred at a call-site as the dot in front of this argument group marks it as *implicit*) and computes a new matrix by applying `f` in pairs to all elements of both inputs. Axioms are declarations of entities that do not have an implementation in THORIN per se. Instead, DSL designers define a *plugin* that consists of an interface (see above) and a C++ implementation that provides domain-specific program transformations and code generators to lower the axioms into the language of lower-level THORIN plugins or to generate target code directly.

IRs like LLVM are not flexible enough to express types like `%matrix.Mat` at all. In MLIR, the designers of plugins (there called *dialects*) need to manually implement the type system in C++ for each dialect operation. But even this has many severe limitations as MLIR does neither support higher-order functions nor polymorphism. For example, `vector.reduction`<sup>1</sup> from the MLIR *vector* dialect hard-codes a set of predefined reduction operations. In THORIN we can define a reduce function that works on any array and appropriate function. For example, here we create a matrix addition with `%matrix.zip elem_add` that we use to reduce an `array_of_matrices` to a single `res_matrix`:

```
.let res_mat = reduce (%matrix.zip elem_add) (initial, array_of_matrices);
```

Note that `%matrix.Mat`’s type is an ordinary function type and, yet, it is a type operator. Furthermore, note that the types of `%matrix.zip`’s arguments `a` and `b` depend on the type variable `T` (polymorphism) and the term variables `n` and `m` which makes `%matrix.Mat (n, m, T)` a *dependent type*. This is possible because at its core THORIN is a minimal, typed  $\lambda$ -calculus that is based upon the  $\lambda$ -cube [7], the Calculus of Constructions (CC) [16], and pure type system (PTS) [37]. This makes THORIN not only very expressive but also simplifies the compiler design in many ways because THORIN only knows a single syntactic category: expressions. However, THORIN’s surface language supports syntactic sugar that THORIN immediately translates into its minimal core language during parsing.

## 1.1 Contributions

In summary, this paper makes the following contributions:

- We introduce the intermediate representation THORIN that allows, through its plugin architecture, to modularly extend the THORIN compiler with custom, domain-specific operations, types, and type operators, as well as domain-specific code transformations—in particular so-called *normalizations* that are eagerly applied by THORIN (Section 2).
- THORIN’s expressive type system is rooted in the Calculus of Constructions and, hence, features dependent-types. These are an overlapping feature of type theory and type systems and, thus, are mainly used in proof assistants and associated with writing complex proof terms. THORIN, however, features dependent types without the hassle of writing complex proof expressions. This is possible because THORIN’s type checker tightly interacts with normalizations and a partial evaluator. THORIN is also a classic compiler in the sense that

<sup>1</sup>see <https://mlir.llvm.org/docs/Dialects/Vector/#vectorreduction-vectorreductionop>

it is built around a two-stage compilation/execution model (compile time & run time). In most other dependently typed languages this line is blurry (Section 3).

- Code transformations are implemented in C++ by manipulating THORIN’s program graph. This is a “sea of nodes”-style [12] IR for a higher-order, PTS-style language. As THORIN only knows expressions, the program graph is also very simple: Each node represents an expression and has outgoing edges to each of its operands and the type it inhabits—which is again an expression. Besides an internal hash set that hash-conses all nodes, THORIN does *not* need any other auxiliary data structures such as instruction lists, basic blocks, CFGs, or special regions (Section 4).
- THORIN performs type checking, inference, and normalization on-the-fly *during program construction*. This way, a plugin/compiler developer can resort to type inference, normalization, and (partial) evaluation upon creating an expression. In addition, THORIN will immediately report any typing errors (Section 5).
- Section 6 presents three case studies that show THORIN’s versatility and ability to host DSLs and generate high-performance code. These include a set of low-level, LLVM-like plugins, that we show perform just as well as if directly using LLVM, a plugin for matching regular expressions (RegExes) which outperforms other popular RegEx engines, and a plugin for automatic differentiation with state-of-the-art performance at a tenth of the complexity.

The THORIN version presented in this paper is actually THORIN 2 and a complete overhaul of its predecessor. See Section 7 for a discussion.

## 2 OVERVIEW

The core of THORIN is a minimal, typed  $\lambda$ -calculus with strong semantics based on the  $\lambda$ -cube, PTS, and CC. We describe the core calculus’ syntax and semantics in more detail in Section 3. What makes THORIN an attractive target for DSLs is its extensibility through plugins. These plugins can define intrinsic operations (in THORIN called “axioms”) on various abstraction levels. Furthermore, plugins can provide transformations and lowering passes to perform (domain-specific) optimizations on each abstraction level and to convert between the levels.

There are several ways how a DSL or a general-purpose language can target THORIN. The most important ones are: an embedded DSL or a DSL compiler may use THORIN’s API to construct the IR. Alternatively, it can communicate with THORIN textually through its surface language. THORIN itself comes with a backend that emits textual LLVM IR. This IR can be further processed by LLVM tools to obtain an executable. Thereby, THORIN currently targets any CPU supported by LLVM.

Since plugins are essential to THORIN’s design, we will use the example of a plugin for matching RegExes to present THORIN’s components.

### 2.1 Plugin Architecture

A plugin consists of two parts:

- (1) A \*.thorin file that contains declarations in THORIN’s surface language of what the plugin exports.
- (2) Code transformations implemented in C++ (see Figure 1).

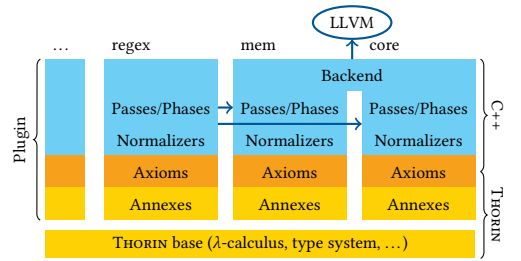


Fig. 1. THORIN architecture overview

```

148 1 .let Char = %core.I8;
149 2 .lam Str (n: .Nat): * = %mem.Ptr «n; Char»;
150 3 .lam Res (n: .Nat): * = [%mem.M, .Bool, .Idx n];
151 4 .let RegEx = Π.[n: .Nat][%mem.M, Str n, .Idx n] → Res n;
152 5
153 6 .ax %regex.any: RegEx;
154 7 .ax %regex.conj: Π.[i: .Nat][«i; RegEx»] → RegEx, normalize_conj, 2;
155 8 .ax %regex.disj: Π.[i: .Nat][«i; RegEx»] → RegEx, normalize_disj, 2;
156 9 .ax %regex.range: «2; Char» → RegEx, normalize_range, 1;
157 10 .ax %regex.not: RegEx → RegEx, normalize_not, 1;
158 11 .ax %regex.quant(optional, star, plus): RegEx → RegEx, normalize_quant, 1;
159 12
160 13 .lam %regex.lit(val: Char) = %regex.range (val, val);
161 14 .let %regex.cls.d = %regex.range ('0', '9'); // similar: %regex.cls.w, %regex.cls.W,
162 15 .let %regex.cls.D = %regex.not %regex.cls.d; // %regex.cls.s, %regex.cls.S
163 16
164 17 .ax %regex.lower_regex: %compile.Pass;

```

Listing 1. RegEx plugin declaration file regex.thorin

```

161 if (auto star_outer = match(regex::quant::star, r))
162   if (auto star_inner = match(regex::quant::star, star_outer->arg())) return star_inner;

```

Listing 2. C++ code that matches  $r^{**}$  and yields  $r^{*}$ 

```

163 world.call<regex::range>(Defs{a, b}, Defs{mem, str, pos})

```

Listing 3. C++ code that constructs a curried call `%regex.range (a, b) (mem, str, pos)`

The `%regex` plugin declares its operations in `regex.thorin` (Listing 1). Keywords are prefixed with a dot to avoid name clashes with ordinary identifiers. All names prefixed with `%regex` are called *annexes*. Before building the plugin’s C++ sources, THORIN *bootstraps* the plugin by parsing `regex.thorin` and generating a C++ header file that declares all annexes as C++ `enums`. This process does *not* involve a “compilation” of the annexes to C++ in any way. The purpose of this is that the C++ part of the plugin can reference an annex via a C++ name instead of a string (see below). Then, the build system compiles the plugin’s sources with the help of the generated header into a shared object `libthorin_regex.so`. Now, other THORIN code can use the plugin:

```

174 .plugin regex; // parse "regex.thorin" and load "libthorin_regex.so"
175 .let pattern = %regex.conj (%regex.quant.plus %regex.cls.w, // '\w+\.[a-z]+'
176                           %regex.lit '.', %regex.quant.plus (%regex.range ('a', 'z')));

```

The `.plugin` `regex` directive instructs THORIN to parse `regex.thorin`. In doing so, THORIN will know the names of all `%regex` annexes plus their types. In addition, THORIN will dynamically load `libthorin_regex.so` to incorporate the plugin’s code transformations into the THORIN compiler. Now that the plugin has been loaded, the `.let`-expression binds the variable `pattern` to a THORIN expression that represents the `RegEx ^\w+\.[a-z]+$` (which matches simple top-level domains).

## 2.2 Plugin Declaration

The `%regex` plugin declares the so-called *axioms* `%regex.conj`, `%regex.disj`, etc. Line 11 declares the `RegEx` quantifiers `%regex.quant.optional`, `%regex.quant.star`, `%regex.quant.plus`. In addition to axioms, plugins can provide supplementary definitions (variables, functions, ...) that are entirely defined in THORIN’s surface language. For example, the digit character class `%regex.cls.d` is just `.let`-bound to `%regex.range ('0', '9')`. Axioms that inhabit a function type, are usable like ordinary functions. However, axioms do *not* provide an implementation in THORIN per se. Their purpose is to denote domain-specific language constructs that the plugin’s code transformations refine. For example, the C++ code in Listing 2 matches

```

192 %regex.quant.star (%regex.quant.star regex)

```

and peels off the superfluous outer quantifier. Note that `regex::quant::star` stems from the auto-generated header and references `%regex.quant.star` from `regex.thorin`. Similarly, plugin developers can access or call annexes (Listing 3).

## 2.3 Types

One of THORIN’s most prominent features is that it does *not* have a special syntactic category for types. Instead, types are also expressions. This has several advantages as we will outline in the following.

First, THORIN does not need special constructs to declare type aliases or type-level functions. An ordinary `.let`-binding (line 1 in Listing 1) defines the alias `Char` for the 8-bit-wide integer type `%core.I8` from the `%core` plugin (Section 6.1.1). `Str` is also a normal function. It expects a size `n` of type `.Nat` and returns `*`: the type of all types. Thus, the type of `Str` is `.Nat → *`. Here, `Str` yields a pointer to an array of size `n` and element type `Char`. The pointer type constructor is an axiom from the `%mem` plugin: `.ax %mem.Ptr: * → *` (Section 6.1.3). The expression `«en; T»` introduces an array *type* while both `en` and `T` are again expressions. We just use the metavariable `T` to suggest that this is a type expression. In particular, the size `en` can be an arbitrarily complex expression and may not necessarily be a compile-time constant as opposed to `n` in C++’s `std::array<T, n>`. For this reason, `«en; T»` is called a *dependent type*, as the type depends on a value. A common misconception is to think of `«en; T»` as a pair consisting of the size `en` and the “actual” array. This is *not* the case. It is just an array which “tracks” its size `en`. Dependent types are mostly known from theorem provers such as Coq or Lean. However, in THORIN we are not concerned about proofs. We are using dependent types to abstract from the size of integer operations or arrays which in turn allows for type-safe variadic functions and polymorphism over array rank while tracking this dependency in the type system (Section 6).

The expression `(e0, ..., en-1)` forms a tuple *value* while `[T0, ..., Tn-1]` forms a tuple *type*. Hence, the function `Res` returns a tuple *type* which models the result type of a `RegEx` match. The first element type `%mem.M` stems again from the `%mem` plugin and abstracts from the machine state; any operation that potentially has a side-effect such as memory accesses consumes a machine state, i.e. a value of type `%mem.M`, and produces a new one. This is similar to the `IO` monad in Haskell. The second element type `.Bool` indicates the success or failure of a match. The last element type `.Idx n` is an integer within the range `0n, ..., (n-1)n` and keeps track of the current position within the string to match. Note that this index type guarantees to access the string in bounds.

The type `RegEx` of a `RegEx` matcher is a *dependent* function type. THORIN’s surface language allows for convenient specification of curried function types with complex dependencies via `Π d ... d → T`: Each domain `d` constitutes a curried domain and may as well as the final codomain `T` depend on the *value* of the preceding domains. Here, the second domain and the codomain `Res n` depend on `n` which is a *value* of the first domain `.Nat`. This dependency makes the first domain deducible when calling a function of type `RegEx`. This is why there is a dot in front of the first domain which marks it as *implicit*: THORIN will automatically infer this argument when calling a function of type `RegEx`. The second domain constitutes the “actual” parameters of the matcher: a machine state, a string, and the current position within this string.

The `RegEx` constructors yield matchers of type `RegEx` while potentially composing other matchers. For example, `%regex.not` expects another matcher to negate and `%regex.range` a range given by two `Chars` as showcased in Listing 3: It creates a range pattern and matches it on a string. Note that the size argument `n` is implicit and inferred in both THORIN’s textual representation as well as the C++ code. The junctions `%regex.conj` and `%regex.disj` demonstrate how dependent arrays allow for variadic functions as they expect *i*-many `RegEx` matchers as inputs.

Declaring these constructors with THORIN types automatically gives them the full power of THORIN’s type system. The `%regex` plugin does not have to provide any kind of additional validation functions as would be necessary in MLIR. It also means that the `Str` type constructor defined in `regex.thorin` is not an opaque type that other plugins would need to know about. On the

contrary, THORIN ensures that any argument is valid automatically if and only if it is assignable to the resulting pointer type.

## 2.4 Normalization

Whenever THORIN creates an expression, it is immediately *normalized* (Section 3.1.1). For example, the tuple extraction  $(0, 1, 2)_{\#2_3}$  is right away resolved to 2. In addition, THORIN consistently removes 1-tuples and 1-tuple types. Hence, it does not matter whether a function is specified as `.lam Str(n: .Nat)` or `.lam Str n: .Nat` and whether it is invoked with `Str n` or `Str (n)`.

Normalizations are the backbone of THORIN’s optimizer but also influence type checking. THORIN handles tuples and arrays in a uniform way by normalizing a tuple type `[.Nat, .Nat]` to an array `«2; .Nat»`. This is why THORIN does not need different introduction and elimination constructs for tuples and arrays: The type of  $(0, 1)$  is `«2; .Nat»` and of  $(0, .ff)$  is `[.Nat, .Bool]`; extraction `e#ei` works regardless of whether `e` is an array or a tuple (Section 3.1.3). Thus, `%regex.conj (re1, re2, re3)` types fine: THORIN infers 3 for `i`, types the argument as `«3; RegEx»`, and removes superfluous 1-tuple types from `%regex.conj`’s domains. Additionally, less syntax also implies fewer patterns to match when writing program analyses.

Finally, all expressions are *hash-consed*: Whenever an expression is created, THORIN first checks whether a syntactically equal expression already exists. If this is the case, THORIN will reuse this existing expression. This has the effect that in the C++ implementation two pointers to THORIN expressions enjoy pointer equality if they are syntactically equal.

Axioms can provide their own *normalizers*: local transformations that are considered generally useful. The `%regex` plugin, for instance, merges quantifiers—`r*?`, `r?*?`, `r+?`, `r?+`, `r++`, `r++` all normalize to `r*`—and removes idempotence—`r??` results in `r?`, ditto for `*` and `+`. As the axiom declaration indicates (line 11), the C++ function `normalize_quant`, which is part of `libthorin_regex.so`, is the *normalizer* of this axiom and implements this logic. The actual implementation consists of ~20 lines of C++ code and involves a few matches and building new calls similar to Listing 2 and 3. The `%regex` plugin implements similar normalizations for the other axioms to compute a (pseudo) normal form for RegExes.

By default, THORIN fires the specified normalizer when the last curried argument is applied to an axiom and this is in most cases the desired behavior. However, plugin designers can override this behavior and specify when exactly normalization should happen. The `%regex` plugin is an exception to the default behavior as it wants to normalize, for example, a quantifier as soon as its matcher is applied: `%regex.quant.star regex`. If we waited until all curried arguments were passed, we would entirely miss normalization in some instances or be too late for others:

```
%regex.quant.star (%regex.quant.star regex /*miss*/) (mem, str, pos) /*too late*/
```

To this end, the axiom demands normalization when the first argument is passed to the axiom (last part of line 11). The same counting mechanism applies when matching curried axiom calls in C++ and is the reason why Listing 2 works as intended. Other plugins implement constant folding and various peephole optimizations such as  $x + 0 \triangleright x$  via normalizers.

## 2.5 Lowering

Since axioms are opaque entities without implementation, plugins must somehow provide an implementation. The `%mem` plugin contains an LLVM backend that additionally understands the `%core` plugin (for integer operations) and `%math` plugin (for floating-point operations). Unless other plugins want to ship their own or extend the existing LLVM backend, they need a phase that substitutes axioms unknown to the LLVM backend to low-level code known to the backend (i.e., only using operations from `%mem`, `%core`, and `%math`). What is more, many plugins want to apply



domain-specific transformations on the code in addition to normalizations before lowering its domain-specific axioms. To this end, THORIN provides a sophisticated optimizer and a flexible, modular pass manager. Even compilation phases are exposed as axioms with the help of the `%compile` plugin. This allows users to compose their own compilation pipeline as a THORIN program. However, the details of the optimizer are beyond the scope of this paper.

As outlined above, the `%regex` plugin applies various normalizations to given RegExes. The `%regex.lower_regex` pass written in C++ constructs a nondeterministic finite automaton (NFA) from a pattern, makes it deterministic, and minimizes it. Finally, the pass will generate low-level code that implements the minimized deterministic finite automaton (DFA) and replaces the axiom calls with it. This code is amenable to the LLVM backend. We discuss the performance of the plugin in Section 6.2.

### 3 SEMANTICS

In the following, we first present the formal definition of THORIN including its syntax, static semantics, and normalization rules (Figure 2). In order to keep this presentation as concise as possible, we leave out a few details and full recursion which we will discuss afterward. Then, we introduce THORIN’s partial evaluator that tightly interacts with type checking and normalization.

THORIN’s surface language supports syntactic sugar, which THORIN translates into the core syntax directly during parsing (Figure 3). While we do not explicitly present THORIN’s dynamic semantics, it is straightforward to infer it using the normalization rules as a blueprint.

#### 3.1 Syntax, Typing & Normalization

*Preliminaries.* Since THORIN is a variant of CC, it uses the same syntax for terms, types, and kinds. However, we usually use the metavariable  $e$  to evoke a term expression,  $T$  or  $U$  to evoke a type expression, and  $s$  to evoke a kind expression. For the purpose of representation, we restrict THORIN’s type universes to  $*$  (the type of types) and  $\square$  (the type of all type operators) which are organized in a predicative way (Rules `BASE/IND/BOX`). This means that the relation  $\rightsquigarrow$  computes the “maximum” of  $*$  and  $\square$ . For example,  $*** \rightsquigarrow *$  and  $*\square* \rightsquigarrow \square$ . This avoids well-known paradoxes (like Girard’s paradox) associated with self-referential definitions. THORIN’s actual implementation, however, uses a stratified, countably infinite hierarchy of sorts.

THORIN knows several *binders*. These are expressions that introduce a variable of the form  $x : e$ . We follow *Barendregt’s convention*: No variable is both free and bound; every bound variable is bound *exactly once*. While THORIN’s surface language supports lexical scoping, THORIN’s graph representation (Section 4) actually ensures Barendregt’s convention.

We call a binder *parametric* if the introduced variable occurs free in any subsequent expression of the binder. Otherwise, we call the binder *non-parametric*. There is syntactic sugar for non-parametric binders available which omits the variable altogether (Figure 3).

Substitution  $e[e_b/e_a]$  where  $e_a$  is recursively replaced with  $e_b$  within  $e$  is defined in the usual manner with one addition: *All expressions created along the way are also normalized.*

**3.1.1 Normalization.** Whenever THORIN builds an expression, it will immediately *normalize* it according to  $\underline{e} \triangleright e$ . We *underline* an expression  $\underline{e}$  to denote that it *might not* be normalized. A *non-underlined* expression  $e$  denotes that it is already *normalized*. Non-normalized expressions only exist in the unparsed surface language of THORIN. For this reason, the normalization rules do not include premises to normalize subexpressions. All subexpressions *are* already normalized.

*Example 3.1.* THORIN will never need to normalize  $((e))$ . As soon as  $(e)$  is parsed, `RULE N-TUP1` fires and yields  $e$ . Then, the outer parentheses form  $(e)$  again which will again fire `N-TUP1` and yield  $e$ .

344	$\Gamma ::= \cdot \mid \Gamma, x : T$	Typing Environment	$\boxed{s^* \rightsquigarrow s}$	BASE $\frac{}{\rightsquigarrow *}$
345	$e, \underline{U}, \underline{s} ::= * \mid \square \mid \perp \mid \text{.Nat} \mid \text{.Idx} \mid n \mid i_n$	Star / Box / Bottom / Nat / Idx / Literal		
346	$i, n \in \mathbb{N} \mid \Pi x : \underline{T} \rightarrow \underline{U} \mid \lambda x : \underline{T} @ (e) : \underline{U} = e \mid e \mid e$	Var / Let / Axiom	IND $\frac{\bar{s} \rightsquigarrow s'}{*s \rightsquigarrow s'}$	BOX $\frac{\bar{s} \rightsquigarrow s'}{\square s \rightsquigarrow \square}$
347	$i < n \mid \lfloor x : \underline{T} \rfloor \mid \langle e \rangle \mid \langle x : e; \underline{T} \rangle \mid \langle x : e; e \rangle \mid e \# e$	Pi / Lam / App	Sigma / Tuple / Array / Pack / Extract	
348	$\boxed{\Gamma \vdash e : T}$	STAR $\frac{}{\Gamma \vdash * : \square}$	BOT $\frac{}{\Gamma \vdash \perp : *}$	NAT $\frac{}{\Gamma \vdash \text{.Nat} : *}$
349		IDX $\frac{}{\Gamma \vdash \text{.Idx} : \text{.Nat} \rightarrow *}$	LIT-N $\frac{}{\Gamma \vdash n : \text{.Nat}}$	
350	LIT-I $\frac{i < n}{\Gamma \vdash i_n : \text{.Idx } n}$	VAR $\frac{x : T \in \Gamma}{\Gamma \vdash x : T}$	AX $\frac{\Gamma \vdash T : s \quad \Gamma, x : T \vdash e : U}{\Gamma \vdash \text{.ax } x : T; e : U}$	PI $\frac{\Gamma \vdash T : s_T \quad \Gamma, x : T \vdash U : s_U \quad s_T \ s_U \rightsquigarrow s}{\Gamma \vdash \Pi x : T \rightarrow U : s}$
351		LAM $\frac{\Gamma, x : T \vdash e_f : \text{.Bool} \quad \Gamma, x : T \vdash U \leftarrow e}{\Gamma \vdash \lambda x : T @ (e_f) : U = e : \Pi x : T \rightarrow U}$	APP $\frac{\Gamma \vdash e : \Pi x : T \rightarrow U \quad \Gamma \vdash T \leftarrow e_T}{\Gamma \vdash e \ e_T : U[e_T/x]}$	
352		SIG $\frac{\Gamma \vdash T_0 : s_0 \quad \dots \quad \Gamma, x_0 : T_0, \dots, x_{n-1} : T_{n-1} \vdash T_{n-1} : s_{n-1}}{S_0 \dots S_{n-1} \rightsquigarrow S \quad \Gamma \vdash [x_0 : T_0, \dots, x_{n-1} : T_{n-1}] : s}$	TUP $\frac{\Gamma \vdash e_0 : T_0 \quad \dots \quad \Gamma \vdash e_{n-1} : T_{n-1} \quad [T_0, \dots, T_{n-1}] \triangleright T}{\Gamma \vdash (e_0, \dots, e_{n-1}) : T} \quad n \geq 0$	
353		ARR $\frac{\Gamma \vdash e_n : \text{.Nat} \quad \Gamma, x : \text{.Idx } e_n \vdash T : s}{\Gamma \vdash \langle x : e_n; T \rangle : s'}$	PACK $\frac{\Gamma \vdash e_n : \text{.Nat} \quad \Gamma, x : \text{.Idx } e_n \vdash e : T \quad \langle x : e_n; T \rangle \triangleright U}{\Gamma \vdash \langle x : e_n; e \rangle : U}$	
354		EX-S $\frac{\Gamma \vdash e : [x_0 : T_0, \dots, x_{n-1} : T_{n-1}] \quad T_j' = T_j[e\#0_n/x_0] \dots [e\#(j-1)_n/x_{j-1}] \quad (T_0', \dots, T_{n-1}') \triangleright T \quad \Gamma \vdash T : s \quad \Gamma \vdash e_i : \text{.Idx } n \quad T\#e_i \triangleright U}{\Gamma \vdash e \# e_i : U}$	EX-A $\frac{\Gamma \vdash e : \langle x : e_n; T \rangle \quad \Gamma \vdash e_i : \text{.Idx } e_n}{\Gamma \vdash e \# e_i : T[e_i/x]}$	
355		A-T $\frac{\Gamma \vdash e : T}{\Gamma \vdash T \leftarrow e}$	A-TUP $\frac{\Gamma \vdash T_0 \leftarrow e\#0_n \quad \forall 1 \leq i < n. \Gamma \vdash T_i[e\#0_n/x_0] \dots [e\#(i-1)_n/x_{i-1}] \leftarrow e\#i_n}{\Gamma \vdash [x_0 : T_0, \dots, x_{n-1} : T_{n-1}] \leftarrow e}$	
356	$\boxed{e \triangleright e}$	N-ID $\frac{\text{if no other rule applies}}{e \triangleright e}$	N-LET $\frac{}{\text{.let } x = e; e' \triangleright e'[e/x]}$	N-EX1 $\frac{}{e\#0_1 \triangleright e}$
357		N-TUP1 $\frac{}{(e) \triangleright e}$		
358		N-SIG1 $\frac{}{[e] \triangleright e}$	N-TUP $\beta$ $\frac{}{(e_0, \dots, e_{n-1})\#i_n \triangleright e_i}$	N-PACK $\beta$ $\frac{}{\langle n; e \rangle \# e_i \triangleright e}$
359		N-TUP $\eta$ $\frac{}{(e\#0_n, \dots, e\#(n-1)_n) \triangleright e}$		
360		N-PACKTUP $\frac{n > 1}{(e, \dots, e) \triangleright \langle n; e \rangle}$	N-ARRSIG $\frac{n > 1}{[T, \dots, T] \triangleright \langle n; T \rangle}$	N- $\beta$ $\frac{e_f[e_a/x] \equiv \text{.tt}}{(\lambda x : T @ (e_f) : U = e_b) \ e_a \triangleright e_b[e_a/x]}$
361		N-TUPPACK $\frac{n \in \mathbb{N} \quad x \in FV\{e\}}{\langle x : n; e \rangle \triangleright (e[\theta_n/x], \dots, e[(n-1)_n/x])}$	N-SIGARR $\frac{n \in \mathbb{N} \quad x \in FV\{e\}}{\langle x : n; T \rangle \triangleright [T[\theta_n/x], \dots, T[(n-1)_n/x]]}$	

Fig. 2. Syntax, Typing & Normalization.  $\underline{e}$  might not be normalized;  $e$  is normalized.

376	Bool	non-parametric binders	function/continuation type
377	$\text{.Bool} := \text{.Idx } 2$	$[...] := [...] : T, \dots]$	$T \rightarrow U := \Pi \_ : T \rightarrow U$
378	$\text{.ff} := 0_2$	$\langle e_n; T \rangle := \_ : e_n; T$	$\text{.Cn } T := T \rightarrow \perp$
379	$\text{.tt} := 1_2$	$\langle e_n; e \rangle := \_ : e_n; e$	$\text{.Fn } T \rightarrow U := \text{.Cn } [T, \text{.Cn } U]$
380	anonymous function/continuation		named function/continuation
381	$\lambda x : T : U = e := \lambda x : T @ (\text{.tt}) : U = e$		$\text{.lam } f x : T : U = e := \lambda x : T : U = e$
382	$\text{.cn } x : T := e := \lambda x : T @ (\text{.ff}) : \perp = e$		$\text{.con } f x : T := e := \text{.let } f = \text{.cn } x : T := e$
383	$\text{.fn } x : T : U = e := \text{.cn } (x : T, \text{return} : \text{.Cn } U) = e$		$\text{.fun } f x : T : U = e := \text{.let } f = \text{.fn } x : T : U = e$
384	(a) Simple syntactic sugar		
385	(b) Curried functions/continuations		
386	(c) Curried function/continuation types		

Fig. 3. Syntactic sugar (excerpt). All functions in 3b are equivalent. The type can be expressed by any expression in 3c—they are equivalent, too. The last respective item depicts the completely desugared version. Similar sugar is available for  $\text{.lam}$ ,  $\text{.con}$ ,  $\text{.fun}$ . In addition,  $\text{.lam}$ / $\text{.con}$ / $\text{.fun}$  allow for recursion (Section 3.2).



*Example 3.2.* Consider `.let x = 3; x`. **Rule N-LET** will immediately normalize this expression to 3. In fact, `.let`-expressions *only* exist in the unparsed surface language of THORIN as **N-LET** will eliminate them outright. For this reason, there is no typing rule for a `.let`-expression. Note that all expressions are hash-consed (Section 4) and, thus, appear exactly once in the program graph. In fact, the implementation does not really perform a substitution here. The front-end simply memorizes that `x` refers to a certain subgraph and all uses of `x` will be wired to that subgraph.

Plugins have the opportunity to extend normalizations (see Section 2.4) by adding rules of the form:

$$\%my.axiom\ e_{arg_1} \cdots e_{arg_n} \triangleright e_{something}.$$

Typical examples include constant folding or various identities like  $x + 0 \triangleright x$ .

Normalization rules play not only an important role in THORIN's optimizer but also in its type checker. First, types themselves are normalized. Second, a normalized expression may appear as argument to another type constructor. In particular in the case of dependent types, checking for type equality requires checking for program equivalence (see also Section 3.3).

*Example 3.3.* Consider the function:

```
λ (a: «%core.nat.add (0, n); T»): U = body
```

The `%core` plugin (Section 6.1) normalizes the addition and, hence, simplifies the function to:

```
λ (a: «n; T»): U = body
```

This allows a caller to pass a value of type `«n; T»` to this function. This would be ill-typed without normalization.

Most dependently typed languages suffer from artifacts such as  $f\ (n + 0)$  types in some context but  $f\ (0 + n)$  does not. THORIN's extensible normalization framework is able to mitigate such issues.

*Example 3.4.* The `%core` plugin normalizes both

$$f\ (\%core.nat.add\ (0, n))\ \text{ and }\ f\ (\%core.nat.add\ (n, 0))\ \text{ to }\ f\ n.$$

Normalizations must be deterministic and cycle-free. Cyclic rules such as  $x + x \triangleright 2 \cdot x$  and  $2 \cdot x \triangleright x + x$  can cause THORIN to diverge as it may endlessly oscillate between these two rules. Determinism is achieved automatically since rules are implemented in C++ which is executed deterministically. However, in the future, we would like to permit the plugin designer to directly specify rewrite rules in THORIN. This would allow plugin designers to specify nondeterministic rules but THORIN could also issue warnings or errors in the case of potentially nondeterministic or cyclic rules.

**3.1.2 Nat & Idx.** THORIN has a builtin type `.Nat` inhabited by  $0, 1, 2, \dots$ . Given  $e$  of type `.Nat`, the type `.Idx`  $e$  represents integers within the range  $0_e, \dots, (e-1)_e$  (**NAT/IDX/LIT-N/LIT-I**). For example, type `.Idx 3` has three inhabitants:  $0_3, 1_3, 2_3$ . For convenience, `.Bool` is an alias for `.Idx 2`, as this type has exactly two inhabitants, for whom appropriate aliases are available, too: `.ff` =  $0_2$  (false) and `.tt` =  $1_2$  (true). We use these types for bootstrapping axioms, but they also play an important role within THORIN itself: The *arity*, i.e. the number of elements of a tuple/array, is a value of type `.Nat`, whereas a value of type `.Idx`  $e$  addresses a specific element of a tuple/array with arity  $e$ .

3.1.3 *Tuples, Packs, Arrays &  $\Sigma$ -Types*. Most languages distinguish between array and tuple terms as well as their types. Albeit THORIN does have different syntax, this distinction does not matter much because THORIN normalizes between both representations (see below).

THORIN generalizes dependent pair types to  $n$ -ary dependent tuple types ( $\Sigma$ -types) where the *type* of an element may depend on the *value* of *any preceding* element. A dependent pair—denoted by  $\Sigma x:T.U$  in literature—is written as  $[x: T, U]$  in THORIN. However, THORIN allows for more complex dependent  $\Sigma$ -types:

```
.let Num = [T: *, add: [T, T] → T, mul: [T, T] → T, _0: T, _1: T];
```

The expression  $\langle x: e_n; T \rangle$  forms an *array type* with  $e_n$ -many elements and element type  $T$ . The *arity*  $e_n$  must be of type `.Nat` and may introduce a variable  $x$  whose type is `.Idx  $e_n$`  and may be used inside the *body*  $T$  (`ARR`). **Rule N-SIGARR** compresses *homogeneous* tuple types to arrays whereas **Rule N-ARRSIG** expands *parametric* arrays of *constant* arity to tuple types.

The term  $(e_0, \dots, e_{n-1})$  introduces a *tuple* while  $\langle x: e_n; e \rangle$  introduces a *pack*. Think of a pack as a compressed tuple. Both terms either inhabit a  $\Sigma$ -type or an array. Tuples and packs work analogously to tuple types and arrays but on term level (`TUP/PACK/N-TUPPACK/N-PACKTUP`).

**Example 3.5.** Due to normalization THORIN considers both the non-normalized as well as the normalized expressions as equal:

```
[.Nat, .Nat] ▷ «2; .Nat»
(0, 0) ▷ <2; 0>
```

```
«i:2; F i» ▷ [F 02, F 12]
<i:2; f i> ▷ (f 02, f 12)
```

The term  $e \# e_i$  *extracts* element  $e_i$  from  $e$ . The index  $e_i$  must type as an `.Idx  $e_n$` . If  $e$  types as array with  $e_n$ -many elements, the type of the extract is the body of the array while substituting the array's variable with the given index  $e_i$  (**Ex-A**). If  $e$  types as  $\Sigma$ -type with  $e_n$ -many elements, **Ex-S** resolves the dependencies of the  $\Sigma$ -type by substituting all preceding variables  $x_j$  with  $e \# j_n$ . The type of the extraction is another extract with index  $e_i$  on the tuple formed by the resolved types.

**Example 3.6.** Suppose `nmx` has type  $[n: \text{.Nat}, m: \text{.Nat}, x: F \ n \ m]$ . Then, `nmx\#23` has type  $F \ nm_x \#0_3 \ nm_x \#1_3$ .

Expression	Type
$(0, 1, 2)$	$\langle 3; \text{.Nat} \rangle$
$(0, 1, 2) \# i$	<code>.Nat</code>
$(0, \text{.tt})$	$[\text{.Nat}, \text{.Bool}]$
$(0, \text{.tt}) \# 0_2$	<code>.Nat</code>
$(0, \text{.tt}) \# 1_2$	<code>.Bool</code>
$(0, \text{.tt}) \# i$	$(\text{.Nat}, \text{.Bool}) \# i$
$(\text{.Nat}, \text{.Bool}) \# i$	$\langle 2; * \# i \rangle \triangleright *$
$(0, \text{.Bool}) \# i$	$\downarrow$

Table 1. Typing examples

**Rule N- $\beta$**  eliminates an extract from a tuple with a known index while **N-PACK $\beta$**  eliminates an extract—no matter the index—from a non-parametric pack. **Rule N-TUP $\eta$**  resolves a tuple comprised of a sequence of extracts from the same entity with increasing indices. Finally, THORIN consistently removes 1-tuples (**N-TUP<sub>1</sub>**), 1-tuple types (**N-SIG<sub>1</sub>**), and extracts with  $0_1$  (**N-EX<sub>1</sub>**).

**Example 3.7.** Note that most languages need different syntax to introduce a tuple or an array term such as  $(0, 1, 2)$  vs.  $[0, 1, 2]$  in Rust. THORIN does not need this distinction. As **Table 1** showcases, tuples with homogeneous element types are typed as array. This makes them amendable for extractions with an unknown index. However, tuples with inhomogeneous element types are typed as  $\Sigma$ -type. Extraction with a known index yields the corresponding element type as expected. However, extraction with an unknown index yields a type computation as another extract (third last row). This again yields a type computation as another extract (second last row). An extract with an unknown index on a tuple whose elements do not agree on their sort is ill-typed (last row).

### 3.1.4 Assignable. Consider the pair

`(.Nat, %core.ncmp.1)` of type `[*, [.Nat, .Nat] → .Bool]`.

However, we can also type it as

```
.let Cmp = [T: *, [T, T] → .Bool]
```

which is similar to a trait in Rust or type class in Haskell. Most systems featuring existential or  $\Sigma$ -types require tuples to be *ascribed* such as `(.Nat, %core.ncmp.1):Cmp`. In THORIN, tuples are *not* ascribed. Instead, whenever an expression  $e$  is *assigned* to a variable  $x: T$ , THORIN checks via  $\Gamma \vdash e \leftarrow T$  whether this assignment actually makes sense. This is trivially the case, if  $e$ 's type is  $T$  (**A-T**). Otherwise, **A-TUP** recursively checks whether all elements of a tuple are assignable while successively resolving the dependencies the  $\Sigma$ -type  $T$  may introduce (similar to **Ex-S** discussed above).

*Example 3.8.* In the following code the type checker allows passing `(.Nat, %core.ncmp.1)` to  $f$ , although it expects an instance of `Cmp`: **Rule LAM** asks **Rule A-TUP** whether the given pair is assignable to `Cmp`, which it is in fact.

```
.lam f(T: *, less: [T, T] → .Bool)(x: T): .Bool = less (x, x);
f (.Nat, %core.ncmp.1) 23
```

This behavior gives THORIN a whiff of duck typing since the same tuple is passable to any function to whose domain the tuple is assignable. In addition, this allows for more structure sharing during hash-consing as differently typed tuples would also result in different nodes.

THORIN also provides an `.insert (et, ei, ev)` operation, that we elided in the formal presentation. It non-destructively creates a new tuple where the element at index  $e_i$  has been replaced with  $e_v$ . From a semantic point of view, insertion is a mix of term elimination and introduction as, for instance, `.insert (et, 13, ev)` is the same as `(et#03, ev, et#23)`.

**3.1.5 Functions.** A function  $\lambda x:T@(\epsilon_f):U = e$  has a dependent function type  $\Pi x:T \rightarrow U$  (**LAM**). This means that the *type* of the function's codomain  $U$  may depend on the *value* of the argument  $x$  which in turn must inhabit the domain  $T$ . The callee  $e$  of an application  $e_{\mathcal{T}}$  must type as a  $\Pi$ -type and the given argument  $e_T$  must be *assignable* (see above) to the domain; the type of the application is resolved by substituting  $x$  with the argument in the codomain  $U$  (**APP**). The Boolean term  $@(\epsilon_f)$  is the so-called *filter* and is used for partial evaluation (**Section 3.3**).

*Continuation-passing style (CPS).* Continuations are functions that never return. THORIN models them as functions whose codomain is  $\perp$  (**BOT**)—a type without inhabitants—and provides syntactic sugar (**Figure 3a**) for continuations `(.cn/.con)` and their types `(.Cn)`. Continuations bridge the gap between CFGs and the  $\lambda$ -calculus, as continuations are akin to basic blocks.

*Example 3.9.* Consider the “if diamond” in **Figure 4b**. The expression `(F, T)#cond ()` first selects either the  $F$  or  $T$  continuation. This works because `cond` is of type `.Bool`—an alias for `.Idx 2`. Applying the selected continuation to `()` continues execution there. The  $T$  case invokes  $N_{42}$  while the  $F$  case invokes  $N_{23}$ . In static single assignment (SSA) form,  $f$ ,  $F$ ,  $T$ , and  $N$  would be basic blocks and  $N$  would need a  $\phi$ -function to select either 23 or 42. In CPS the  $\phi$ -function becomes the parameter  $\phi$  of continuation  $N$  while the operands of the  $\phi$ -function become arguments to the appropriate continuation call [5, 26].

Continuations may be used (mutually) recursive (**Section 3.2**) to model arbitrary, unstructured control flow. Note that  $f$  in **Figure 4b** is a continuation and the return point is explicit by passing it to  $f$  as another continuation. This idiom is so common that THORIN introduces `.fn/.fun/.Fn` as sugar (**Figure 3a**). For example, in **Figure 4b** we can write `.fun f(cond: .Bool) → .Nat = /*...*/`,

instead. Finally, THORIN provides syntactic sugar for *curried* functions & continuations (Figure 3b-c). Note that all curried function groups except the last one are in direct style in the case of curried continuations.

### 3.2 Full Recursion

Named functions (Figure 3a) are in fact more powerful than ordinary `.let` bindings, as they allow for (mutual) recursion and, hence, are more like `letrec` in other languages.

```
.lam forever(x: .Nat): .Nat = forever x;
```

This is an extension of the calculus presented so far. Internally, the recursive function `forever` is represented as a cyclic graph (Section 4). THORIN allows recursion in both direct style (Example 3.11) and CPS:

*Example 3.10.* The following loop counts from 0 to 42. If its parameter `i` (the “ $\phi$ -function”) is less than 42, `i` is incremented and loop recurses. Otherwise, it exits.

```
.con loop(i: .Nat) =                                     // i =  $\phi(0, i + 1)$ 
  .con body() = loop (%core.nat.add (i, 1)); // recurse
  .con exit() = foo i;                               // pass last instance of i = 42 to foo
  (exit, body)#(%core.ncmp.l (i, 42)) (); // if i < 42 then body() else exit()
loop 0;                                              // loop entry
```

### 3.3 Partial Evaluation & $\beta$ -Equivalence

The so-called *filter*  $@(e_f)$  of a function is a Boolean expression that may depend on the function’s parameter and is used for partial evaluation ( $N\text{-}\beta$ ). For each call-site, THORIN instantiates the filter by substituting the function’s variable with the call-site’s argument. Remember that substitution also normalizes. If this syntactically yields `.tt`, THORIN will  $\beta$ -reduce this call-site—potentially recursively inlining more function calls. This mechanism allows for compile-time specialization.

Since THORIN has a dependent type system featuring full recursion (Section 3.2), type-checking becomes undecidable in general since checking for  $\beta$ -equivalence is. This is arguably the most concerning issue for picking up dependent types in more mainstream programming languages. THORIN tackles this issue by (partially) evaluating functions according to the filters during normalization.

*Example 3.11.* The following recursive function `pow`, will be recursively  $\beta$ -reduced, if the exponent `b` is a compile-time constant. Note that `f`’s parameter `x` has a dependent type.

```
.lam pow(a b: .Nat)@(%core.pe.known b): .Nat =
  (%core.nat.mul (a, pow(a, %core.nat.sub (b, 1))), 1)#(%core.ncmp.e (b, 0));

.lam f(n: .Nat, x: «%core.nat.mul (n, %core.nat.mul (n, n)); .Nat»): [] = ();
.lam g(m: .Nat, y: «pow (m, 3); .Nat»): [] = f (m, y);
```

As `g` calls `f` with `(m, y)`, THORIN has to check whether `y`’s type is assignable to `g`’s domain ( $P_1$ ). Rule  $N\text{-}\beta$  determines that `pow`’s filter evaluates to `.tt` with the given argument `(m, 3)`. This causes THORIN to immediately  $\beta$ -reduce `pow (m, 3)` which will result in a call-site `pow (m, 2)`. Rule  $N\text{-}\beta$  determines again that the filter yields `.tt`, causing another  $\beta$ -reduction. The type checker now sees `y`’s type as

```
«%core.nat.mul (m, %core.nat.mul (m, m)); .Nat»
```

which is assignable to `g`’s domain. Using a `.ff` filter for `pow` would result in a type error, as «`pow (m, 3); .Nat`» is different from `x`’s type. Using a `.tt` filter for `pow` is dangerous: While the example above would still work, a call-site like `pow (i, j)`, where `j` is not a compile-time constant, would cause THORIN to diverge, as it would endlessly  $\beta$ -reduce new calls to `pow`. However, this termination behavior is not random and completely transparent to the programmer. It depends on the specified filters as well as how they are used—as stated by the programmer.

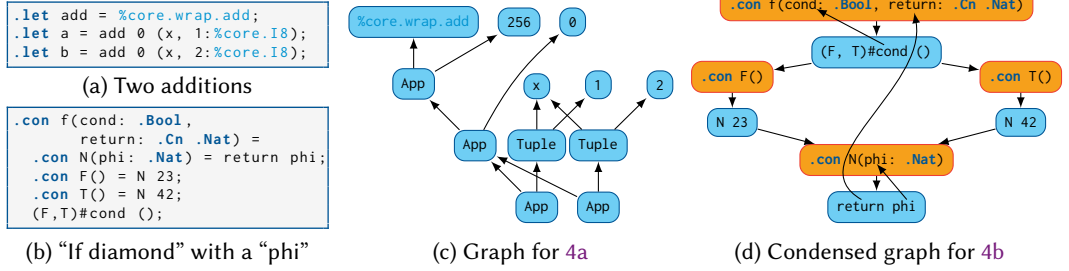


Fig. 4. Thorin’s surface language vs. graph-based representation. Type edges are elided. Immutables are in blue, mutables in orange.

Coupling partial evaluation with normalization that in turn interacts with type checking is a novel approach to dependent type checking. Moreover, partial evaluation filters clearly determine which parts of a program are evaluated at compile time and which ones remain in the compiled program. This distinction is somewhat unclear in many dependently typed languages—Iris 2 is an exception (Section 7).

Note that the rules in Figure 2 do *not* contain a conversion rule that states that two terms are equal modulo  $\beta$ -equivalence. The typing rules are syntax-directed and deterministic similar to Pollack and Poll [47]. The normalization rules unambiguously state where normalization and, in particular,  $\beta$ -reduction happens during type checking.

**3.3.1 Default Filter Policy.** Elided partial evaluation filters default to `.tt` except for the final continuation in a (curried) `.cn/.con/.fn/.fun` function. The filter of this final continuation defaults to `.ff` (Figure 3). This has the effect that the actual computations are deferred to runtime while other abstractions such as type abstractions are specialized at compile time [59]—similar to C++ templates, Rust generics, polymorphism in MLton [51] but unlike Java generics. However, programmers can override the default behavior by explicitly specifying a filter as in Example 3.11.

### 3.4 Type Safety

THORIN is to a large extent based on CC which has a well-established metatheory and, in particular, progress and preservation are known properties [15, 52]. Custom axioms/normalizers, however, may add unsafe features like, for example, `%core.bitcast` or all memory access axioms from `%mem` (see Section 6.1). Note also that many (intermediate) languages like LLVM, MLIR, etc. do not have an official formal description of their semantics at all. A full mechanized formalization of THORIN on top of CC is subject to future work.

## 4 GRAPH REPRESENTATION

THORIN’s implementation is based upon the “sea of nodes” concept [12]. This means that THORIN’s internal representation is a data dependence graph where each node in the graph represents an expression. If  $e_s$  is a subexpression/operand of an expression  $e$ , the graph contains an edge  $e \rightarrow e_s$ . For example, an `App` has two operands—the callee and argument—and, hence, two outgoing edges while a  $\Sigma$ -type with three element types has three operands and, hence, three outgoing edges. In addition, every expression  $e$  has exactly one type  $T$  it inhabits. This relation is modeled with another edge  $e \rightarrow T$  in the program graph. Note once again that  $T$  is an expression. Hence, there is only *one* graph that contains both terms and types as well as any dependencies between them via edges; in particular, we may have types that depend on terms due to dependent types. Furthermore, this graph is *complete*. This means that this graph comprises the whole semantics of a THORIN

program and—apart from an internal hash set to hash-cons all nodes (Section 2.4)—THORIN does *not* rely on any other auxiliary data structures like instruction lists, basic blocks, special regions, or CFGs. `.let`-expressions and even explicit function nesting only exist in the surface language. Thus, a THORIN graph *solely* consists of a large number of nodes that “float” in a complex network, resembling a sea—hence the term “sea of nodes”.

*Example 4.1.* Figure 4c depicts the THORIN graph of Figure 4a. Note that the `.let`-expressions in the surface language are absent in the graph and the curried call `%core.wrap.add /*256*/ 0`, where the inferred implicit argument is now explicit, is shared by both additions. In fact, all 8-bit wide additions with mode 0 (Section 6.1.1) will reuse this subgraph. Similarly, `x` is shared by both argument tuples of the additions. What is not shown in Figure 4c are the type edges. For example, the bottom `Apps` and the literals 1 and 2 point to a subgraph that constitutes `.Idx 256`.

*Immutable vs. Mutable.* So far, we have discussed expressions that we create by first building their operands, and then the actual nodes. THORIN calls these expressions *immutable*: Once constructed, they cannot be changed later on. Hence, immutables form a directed acyclic graph (DAG). In order to allow for recursion, we have to somehow form a cyclic graph. For this reason, there are also *mutable* expressions (Table 2). For mutables, we first build the node and set its operands later on. This enables us to form a cyclic graph and therefore enables recursion. Mutables also allow changing their operands later on—hence their name. This also means that a mutable will not be hash-consed as opposed to immutables. But THORIN will check mutables for  $\alpha$ -equivalence, if necessary. Finally, all expressions (both on term and type level) that introduce variables are also modeled as mutables.

Immutable	Mutable
build operands first, then the actual node	build the node first, then set the operands
operands form a DAG	operands may be cyclic
hash-consed	each new entity is fresh
non-parametric	may be parametric

Table 2. Immutable vs. mutable nodes

*Example 4.2.* Consider the THORIN graph in Figure 4d of Figure 4b. In order to access `f`’s and `N`’s variables, these continuations are mutables. In fact, most of the time we build functions as mutables even though `F` and `T` could be modeled as immutables in this particular case, as they are neither used recursively nor are their variables accessed.

## 5 TYPE CHECKING, INFERENCE, AND NORMALIZATION ON-THE-FLY

THORIN’s implementation performs normalization, type checking, and inference of implicits as soon as a new expression is constructed. This eagerness has the advantage that a compiler/plugin developer will immediately notice if they incorrectly plugged together some expressions that are ill-typed.

### 5.1 $\alpha$ -Equivalence

This means that type checking, inference, and normalization must also work on open terms, i.e., in the presence of free variables. The implementation boils down to checking for  $\alpha$ -equivalence modulo free variables inside of the *assignable* relation. Here we know that *either* two expressions *must* be  $\alpha$ -equivalent *or* there is a type error. For this reason, THORIN optimistically assumes during *type checking* that any free variables are  $\alpha$ -equivalent. Eventually, all terms will be closed where any remaining issues will be found.

*Example 5.1.* During type checking THORIN considers the following expressions as  $\alpha$ -equivalent:

$$\lambda(a: .\text{Nat}): .\text{Nat} = b \quad \text{and} \quad \lambda(x: .\text{Nat}): .\text{Nat} = y.$$



*Example 5.2.* For *normalization*, however, this assumption is unsound. **N-PackTup** cannot simply normalize the following expression as *b* and *y* may be bound differently:

$$(\lambda(a: \text{.Nat}): \text{.Nat} = b, \lambda(x: \text{.Nat}): \text{.Nat} = y) \not\approx \langle 2; \lambda(a: \text{.Nat}): \text{.Nat} = b \rangle$$

As discussed in [Section 4](#), pointer equality of two THORIN expressions implies normalized, syntactic equivalence. Both expressions refer to the same object. However, pointer equality does not necessarily imply  $\alpha$ -equivalence when free variables are involved. For this reason, THORIN only resorts to pointer equality when checking for  $\alpha$ -equivalence in the absence of free variables.

*Example 5.3.* Although the bodies of the following functions enjoy pointer equality, these functions are *not*  $\alpha$ -equivalent as *x* is bound in the first function and free in the second one:

$$\lambda(x: \text{.Nat}): \text{.Nat} = x \quad \text{and} \quad \lambda(y: \text{.Nat}): \text{.Nat} = x$$

## 5.2 Type Inference

Whenever a function with an implicit argument is invoked, the function will first be applied with a placeholder.

*Example 5.4.* Consider the annex `%core.minus` that computes unary minus of its *s*-sized `.Idx` argument *a* ([Figure 5a](#)). Note that *s* is implicit. Furthermore, the operation expects a so-called mode of type `.Nat` (see [Section 6.1](#) for details). As soon as we apply the first explicit argument *mode* to `%core.minus`, THORIN will insert a fresh placeholder—let us say `?5`—as implicit argument in between:

```
%core.minus ?5 mode
```

This yields the type `.Idx ?5 → .Idx ?5` (see [App](#)). When we apply `42256` as third argument, **App** triggers the *assignable* relation. Rules **A-T** and **A-Tup** additionally match placeholders with the provided argument and fill out any gaps. So here, we find: `?5 = 256`:

```
%core.minus /*256*/ mode 42256
```

Due to the implicit `.tt` filter, THORIN will  $\beta$ -reduce the application to:

```
%core.wrap.sub /*256*/ mode (%core.idx 256 mode 0, 42256)
```

This in turn will be normalized to `214256` (two's complement of -42).

## 5.3 C++ Interface

Type checking, inference, and normalization happens *regardless* of whether THORIN is used through its surface language or its C++ interface. In fact, the THORIN parser directly emits a THORIN graph.

*Example 5.5.* To better bring the point across how THORIN behaves when controlling from C++, reconsider [Example 5.4](#). We can construct the same call using the C++-API:

```
const Def* res = world.call<core::minus>(mode, world.lit_idx(256, 42));
```

Now, type checking, inference, and normalization will happen as discussed above. Thus, the C++ variable *res* (of static type `const Def*` and dynamic type `const Lit*`) will point to a literal that represents 244 of type `.Idx 256`. In other words, creating THORIN programs from C++ triggers the usual normalization that may in turn cause a Turing-complete (partial) evaluation of the THORIN program.

*Example 5.6.* When trying to construct an ill-typed THORIN program from C++, THORIN will throw a C++ exception as in the following code that erroneously passes a `.Nat` instead of an `.Idx` literal to `%core.minus`:

```
const Def* res = world.call<core::minus>(mode, world.lit_nat(42)); // %core.minus mode 42
```

## 6 EVALUATION

This section showcases THORIN’s flexibility by discussing some of the plugins we have already developed and evaluating their performance. In particular, we want to show that

- THORIN is able to achieve the same performance as other low-level approaches like C/LLVM.
- THORIN is able to fully remove carefully crafted abstractions.
- designing code analyses/transformations in THORIN is no more complicated than in traditional compiler IRs like LLVM.

If not mentioned otherwise, we ran all tests using a single thread on an AMD Ryzen 7 3700X supported by 64GB of DDR4@2133MT/s RAM. We used LLVM/Clang 15.0.7 for the C sources as well as the LLVM code emitted by THORIN.

### 6.1 Low-Level Plugins

In this section, we discuss the design of three low-level plugins, which are deliberately modeled on LLVM: First of all, LLVM is a well-thought-out low-level IR; second, it allows easy mapping of axioms defined in these plugins to LLVM instructions in THORIN’s LLVM backend. This backend also lives within a plugin. Other more high-level plugins eventually lower their axioms to those of these low-level plugins. The `%core` plugin introduces integer operations, the `%math` plugin a floating-point type operator and operations for it, and the `%mem` plugin side-effects, memory operations, and pointer arithmetic.

**6.1.1 The core Plugin.** This plugin (Figure 5a) defines arithmetic operations (`%core.nat`) and comparisons (`%core.ncmp`) for `.Nat` and integer operations on values of type `.Idx s`. All integer operations abstract over the size `s`. The plugin introduces a set of convenient aliases `%core.i8`, `%core.i16`, etc. for common `.Idx` types. Like in LLVM, it is up to the operation to decide whether a specific `.Idx s`-typed value is signed or unsigned. For example, the right-shift `%core.shr` comes in two flavors: arithmetic (with sign extension), and logical (without sign extension). Integer operations like `%core.wrap.add` take an overflow mode `m`. This mode dictates how overflow is handled (wrap-around or undefined behavior for both signed and unsigned overflow). Along the same lines, `%core` introduces all 4 unary (`%core.bit1`), all 16 binary bitwise (`%core.bit2`), as well as the usual signed/unsigned comparison (`%core.icmp`) operations. The `%core.conv` axiom converts between different-sized signed or unsigned values via truncation, zero-, or sign-extension whereas `%core.bitcast` allows for arbitrary (potentially unsafe) casts.

For convenience, there is a function `%core.minus` available that builds a unary minus by subtracting the given operand `a` from `0` (see also Section 5.2-5.3). This function will always be inlined due to the default `.tt` filter (Section 3.3.1). Note that in LLVM, MLIR, and similar IRs such helpers are usually implemented as C++ code that directly emit the desired code snippet as their type systems cannot express polymorphic/dependently typed functions. The THORIN function `%core.minus` on the other hand can be called and passed around just like any other axiom.

**6.1.2 The math Plugin.** This plugin (Figure 5b) introduces a type operator `%math.F` that expects the number of significant precision bits `p` and exponent bits `e` over which all `%math` operations abstract. Additionally, most operations expect a mode `m` that fine-adjusts how strictly they should obey the IEEE-754 standard for floating-point transformations. For example, you may choose to allow reassociation of floating-point operations or ignore NaNs or infinity. The axioms for the actual operations such as `%math.arithmetic`, `%math.trigonometric`, or floating-point comparisons (`%math.cmp`) and convenience wrappers like `%math.minus` are straightforward.

```

785 .let %core.i8 = 0x100; // similar: i16, i32, i64
786 .let %core.I8 = .Idx %core.i8; // I16, I32, I64
787 // Create literal of type .Idx s from l while obeying mode m
788 .ax %core.idx:  $\Pi[s: \text{.Nat}][m: \text{.Nat}][l: \text{.Nat}] \rightarrow \text{.Idx } s$ ;
789 .ax %core.nat(add, sub, mul):  $[a \ b: \text{.Nat}] \rightarrow \text{.Nat}$ ;
790 .ax %core.ncmp(/*...*/):  $[a \ b: \text{.Nat}] \rightarrow \text{.Bool}$ ;
791 .ax %core.bit1(f, neg, id, t):  $\Pi[s: \text{.Nat}][m: \text{.Nat}][a: \text{.Idx } s] \rightarrow \text{.Idx } s$ ;
792 .ax %core.bit2(/*...*/):  $\Pi[s: \text{.Nat}][m: \text{.Nat}][a \ b: \text{.Idx } s] \rightarrow \text{.Idx } s$ ;
793 .ax %core.wrap(add, sub, mul, shl):  $\Pi[s: \text{.Nat}][m: \text{.Nat}][a \ b: \text{.Idx } s] \rightarrow \text{.Idx } s$ ;
794 .ax %core.shr(a, 1):  $\Pi[s: \text{.Nat}][a \ b: \text{.Idx } s] \rightarrow \text{.Idx } s$ ;
795 .ax %core.icmp(/*...*/):  $\Pi[s: \text{.Nat}][a \ b: \text{.Idx } s] \rightarrow \text{.Bool}$ ;
796 .ax %core.conv(s, u):  $\Pi[ss: \text{.Nat}][ds: \text{.Nat}][\text{.Idx } ss] \rightarrow \text{.Idx } ds$ ;
797 .ax %core.bitcast:  $\Pi[S: *][D: *][S] \rightarrow D$ ;
798 .lam %core.minus .(s: .Nat)(m: .Nat)(a: .Idx s): .Idx s =
799   %core.wrap.sub m (%core.idx s m 0, a);

```

(a) The %core plugin

```

800 .ax %math.F:  $[p \ e: \text{.Nat}] \rightarrow *$ ;
801 .let %math.f64 = (52, 11); // similar: f16, f32, ...
802 .let %math.F64 = %math.F %math.f64; // F16, F32, ...
803 .ax %math.arith(add, sub, mul, div, rem):
804    $\Pi[p \ e: \text{.Nat}][m: \text{.Nat}][a \ b: \text{.Idx } s] \rightarrow \text{.Idx } s$ ;
805 .ax %math.tri(/*...*/):  $\Pi[p \ e: \text{.Nat}][m: \text{.Nat}][a: \text{.Idx } s] \rightarrow \text{.Idx } s$ ;
806 .ax %math.cmp(/*...*/):  $\Pi[p \ e: \text{.Nat}][m: \text{.Nat}][a \ b: \text{.Idx } s] \rightarrow \text{.Bool}$ ;
807 .lam %math.minus .(p \ e: .Nat)(m: .Nat)(a: %math.F (p, e)): %math.F (p, e) =
808   %math.arith.sub m (0: (%math.F pe), a);

```

(b) The %math plugin

```

809 .ax %mem.M: *;
810 .ax %mem.Ptr: *  $\rightarrow$  *;
811 .ax %mem.alloc:  $\Pi[T: *][\text{.Nat}] \rightarrow [\text{.Nat}, \text{.Nat}]$ ;
812 .ax %mem.free:  $\Pi[T: *][\text{.Nat}, \text{.Nat}] \rightarrow \text{.Nat}$ ;
813 .ax %mem.load:  $\Pi[T: *][\text{.Nat}, \text{.Nat}] \rightarrow \text{.Nat}$ ;
814 .ax %mem.store:  $\Pi[T: *][\text{.Nat}, \text{.Nat}, \text{.Nat}] \rightarrow \text{.Nat}$ ;
815 .ax %mem.slot:  $\Pi[T: *][\text{.Nat}, \text{.Nat}] \rightarrow [\text{.Nat}, \text{.Nat}]$ ;
816 .ax %mem.lea:  $\Pi[n: \text{.Nat}, Ts: \langle n; * \rangle][\text{.Nat}, \text{.Nat}] \rightarrow \text{.Nat}$ ;

```

(c) The %mem plugin

```

816 .ax %autodiff.AD: *  $\rightarrow$  *; // marks type-level transformation
817 .ax %autodiff.ad:  $\Pi[T: *][T] \rightarrow \text{.Nat}$ ; // marks term-level transformation

```

(d) The %autodiff plugin

Fig. 5. Selection of plugins we implemented (excerpts). Normalizer information has been elided.

6.1.3 *The mem Plugin.* This plugin (Figure 5c) introduces a type %mem.M to abstract from the machine state (Section 2.3). We expose axioms to allocate memory, load from, and store values into previously allocated pointers, and to index into pointers that point to compound data types. Most axioms expect a machine state and additional arguments like the pointer to operate on, and return a memory instance together with the produced results like the loaded value. The %mem.lea<sup>2</sup> axiom performs pointer arithmetic. Since this instruction does not have side effects as it does not directly interact with memory, no machine state is needed. The %mem.lea axiom takes a pointer to a tuple or array and the offset i in which it wants to index. The result is the pointer to the i-th element.

<sup>2</sup>The name is inspired by the x86 assembly instruction and its semantics is similar to getelementptr in LLVM but arguably more streamlined.

Benchmark	Input size	C [s]	IMPALA [s]
aobench	–	$0.667 \pm 0.014$	$0.659 \pm 0.004$
fannkuch	12	$27.709 \pm 0.210$	$26.301 \pm 0.186$
fasta	25,000,000	$0.711 \pm 0.011$	$0.814 \pm 0.013$
mandelbrot	5,000	$1.393 \pm 0.011$	$1.390 \pm 0.010$
meteor	2,098	$0.036 \pm 0.001$	$0.036 \pm 0.004$
nbody	50,000,000	$3.621 \pm 0.031$	$2.825 \pm 0.021$
pidigits	10,000	$0.371 \pm 0.005$	$0.370 \pm 0.004$
spectral	5,500	$1.387 \pm 0.011$	$1.409 \pm 0.014$
regex	–	$4.101 \pm 0.037$	$4.128 \pm 0.044$
reverse	–	$0.744 \pm 0.025$	$0.714 \pm 0.029$
geom. speedup	–	1	1.020

Table 3. Comparing the original Benchmark Game C version with our IMPALA implementation that uses THORIN.

*Example 6.1.* In the following listing, `p` points to a 3-tuple and `%mem.lea` indexes into the second element (`l3`) resulting in the pointer `q`.

```
.let Foo: * = [%core.I32, %core.I16, %core.I8];
.let p: %mem.Ptr Foo = /*...*/;
.let q: %mem.Ptr %core.I16 = %mem.lea /*(3, (%core.I32, %core.I16, %core.I8))*/ (p, l3);
```

**6.1.4 Evaluation.** In this experiment, we want to confirm that the low-level plugins perform as well as directly using LLVM. To this end, we ported the code generator of the research language IMPALA [33] to THORIN 2. We ported the fastest available C implementations that were neither manually vectorized nor parallelized from *The Computer Language Benchmarks Game* [58] to IMPALA. In addition, we ported the publicly available aobench [21]. Table 3 shows that the performance is as expected nearly identical to the original C versions (with two slight outliers—one for the better, the other one for the worse). Note that the regex benchmark does *not* use the %regex plugin.

## 6.2 Regular Expressions

We have already discussed the %regex plugin in Section 2: The plugin defines a set of axioms representing ranges of literals, consecutive elements (conjunction), alternative elements (disjunction), negation, and quantifiers. These are sufficient to define a useful set of compile-time regular expressions. While we only test this plugin with THORIN’s surface language, one could easily integrate THORIN with this plugin as a code generator into a RegEx parser and either generate code for a RegEx at compile time or just-in-time (JIT)-compile at run time.

This showcase demonstrates that THORIN provides a powerful core whose extensibility indeed makes implementing a DSL with intrinsic, normalizable expressions and domain-specific optimizations very straightforward. Our %regex plugin provides a legalization pass `lower_regex` that is written in C++ and receives the normalized RegEx pattern in its opaque form. This allows the pass to use the exhaustive understanding of the RegEx to generate an optimized matcher using finite

RegEx Engine	LoC	Compile time [ms]	Match time [μs]
CTRE	4,153	1,677	4,736
std::regex	4,874	2,207	10,151
pcre2	85,879	67	3,882
pcre2-jit			1,308
hand-written C	102	45	886
THORIN	965	145	640

Table 4. Comparison of several RegEx engines

<pre> 883 .fun (x: T, y: T): T = 884   .let z = x + y;       return (x * z); </pre>	<pre>       .fun ((x, x*): [T, T → T], (y, y*): [T, T → T]): [T, T → T] = 885       .let z, z* = (x + y, .fn (s: T): T = return (x*(s) + y*(s))); 886       return (x * z, .fn (s: T): T = return (x*(z*s) + z*(x*s))); </pre>
---	--

Fig. 6. Each computation of the original program (left) is augmented with a backpropagator marked with  $\square^*$  (colored in blue, right). The differentiated function returns  $x \cdot (x + y)$ ,  $\lambda s. s \cdot (2x + y, x)$ . THORIN-like pseudocode.

automatons. To do so, the pass first translates the RegEx into a NFA, further converts this to a DFA [63, 2.2], and then minimizes it [25]. Finally, the minimal DFA is translated into low-level control flow in THORIN’s IR that is purely based on integer comparisons as well as jumps between state continuations. The generated code is put into action by the pass as it replaces the pattern application with a call to the newly generated matcher function.

In total, the RegEx plugin only encompasses 943 lines of C++ code and 22 lines of THORIN code. We compare our RegEx implementation in THORIN with Compile Time Regular Expression (CTRE) [18], `std::regex` [20], as well as Perl-compatible Regular Expressions 2 (PCRE2) [24] in an interpreted and a JIT compiled variant. From the lines of code (LoC) numbers in Table 4 we observe that THORIN’s implementation is at least one order of magnitude less complex. Despite the rather low complexity, our engine outperforms state-of-the-art RegEx engines. It is even 28% faster than a manually written, low-level matcher that we implemented in C and matches only the specific pattern below. Most likely, our C implementation contains some redundant checks, but spotting them is very hard in such low-level code that comes from manually writing a complicated matcher.

The listed LoC include actual code lines exclusively.<sup>3</sup> “Compile time” is the execution time of clang++ and for THORIN the THORIN frontend and optimizer. All benchmarks test the following RegEx on 10,215 E-Mail addresses that are matched by this pattern and 450 that are not [48]:

```

^ [a-zA-Z0-9] (?: [a-zA-Z0-9] * [ _ \ - ] * [a-zA-Z0-9] ) * [a-zA-Z0-9] * @ [a-zA-Z0-9] (?: [a-zA-Z0-9] * [ _ \ - ] * [a-zA-Z0-9] ) * [a-zA-Z0-9] * \. (?: (?: [a-zA-Z0-9] * [ _ \ - ] * [a-zA-Z0-9] ) * [a-zA-Z0-9] + \. ) * [a-zA-Z] [a-zA-Z] + $

```

### 6.3 Automatic Differentiation

The `%autodiff` plugin implements reverse-mode automatic differentiation (AD)—a prominent technique to compute the derivatives of code. These derivatives are used to optimize parameters using gradient descent methods in, for example, machine learning frameworks.

A popular way to implement reverse-mode AD in functional languages is to use *backpropagators* [45]. Each function  $f$  is augmented to return a function  $f^*$ , the backpropagator, in addition to its original result (see Figure 6 for an example):

$$Df(x, x^*) := (f(x), \lambda a. x^* (f'(x) \cdot a))$$

By this technique, AD effectively builds a list of backpropagators: The backpropagator of  $f$  uses  $x^*$  which is the backpropagator of the function that was used to compute the argument of  $f$ . In turn,  $f$ ’s backpropagator is passed to all functions that take the value  $f$  returned as an argument. The derivative is then computed by invoking the backpropagator of the end result with 1.

The `%autodiff` plugin expresses this transformation as a set of local rewrite rules that are applied in a bottom-up fashion to transform the original function into its derivative. Besides simple functions, the `%autodiff` plugin handles a rich set of features including non-scalar types, pointers, and higher-order recursive functions. The advantage of the local-rewrite approach to AD is that it is simple to implement and modular in the sense that AD of a compound language element translates to differentiating its constituents.

This modularity allows us to extend the implementation at a high level by simply specifying the derivative of another operation. For example, we have also implemented a `%matrix` plugin (similar

<sup>3</sup>LoC as reported by `cloc`, treating THORIN code as Rust

to the example in Section 1) that specifies derivatives for common matrix operations without tainting the `%autodiff` plugin with dependencies on the `%matrix` plugin.

The use of higher-order functions in the derivative computation makes it more challenging for the compiler to produce efficient code. However, THORIN’s optimizer and partial evaluator are able to remove this overhead entirely in the benchmarks that we considered as our evaluation below shows. In particular, THORIN’s design resolves many practical implementation problems basically for free: Most notably, hash-consing merges identical backpropagator invocations, and partially evaluating backpropagators removes most of the boilerplate that is introduced by the newly introduced nested function calls.

To show that THORIN is able to succinctly optimize the code that results from applying backpropagator-based AD, we compare it against the state-of-the-art AD frameworks PYTORCH [44] and ENZYME [40]. PYTORCH builds dynamic computation graphs via operator overloading in Python. ENZYME differentiates LLVM code during compilation.

*Setup.* We evaluate the frameworks on Microsoft’s ADBENCH suite [53]. We use the Gaussian mixture model (GMM), bounded analysis (BA), and a long short-term memory neural network (LSTM) from the ADBENCH suite. Additionally, we compare the approaches on a network for the MNIST classification task [17].

For a fair comparison to ENZYME, we instruct the `%matrix` plugin to generate low-level, straight-forward loop nests. Alternatively, the plugin could also employ highly tuned BLAS [9] routines, instead. But we are interested in how well the `%autodiff` plugin copes with low-level loop nests as this abstraction layer corresponds to the knowledge ENZYME obtains via LLVM’s scalar evolution analyses.

**6.3.1 Running Time.** We ran the benchmarks using a single thread on an AMD Ryzen 7 5800X with 32 GB of DDR4@2400MT/s RAM. Figure 7 compares the speedup/slowdown of ENZYME/PYTORCH compared to `%autodiff` as baseline. For the THORIN implementations, we write the benchmarks in Impala (Section 6.1.4) that compiles to THORIN and uses the `%autodiff` AD compiler pass. Finally, THORIN emits an LLVM file that we feed to CLANG to generate an executable. The ENZYME implementation<sup>4</sup> is written in C++ and uses the ENZYME LLVM pass. The PYTORCH

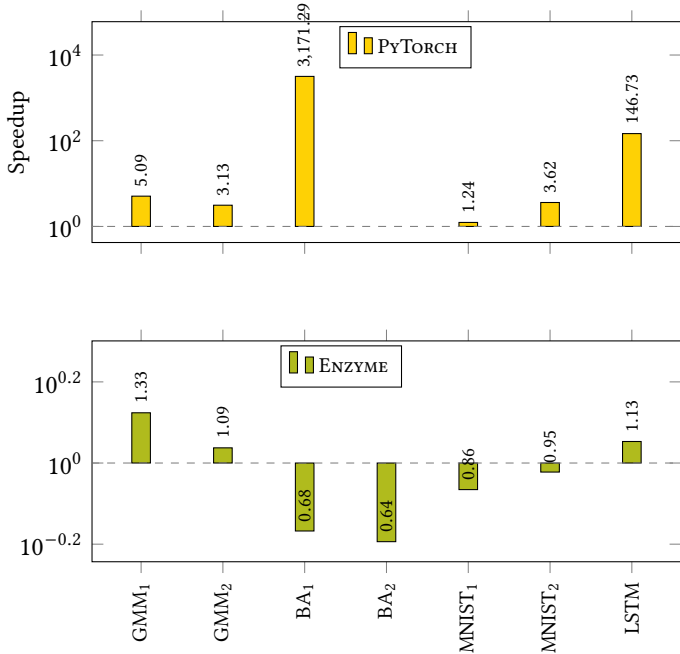


Fig. 7. Speedup  $t/t_{\text{thorin}}$  of THORIN vs. ENZYME, and PYTORCH

<sup>4</sup><https://github.com/EnzymeAD/Enzyme/tree/main/enzyme/benchmarks/ReverseMode>



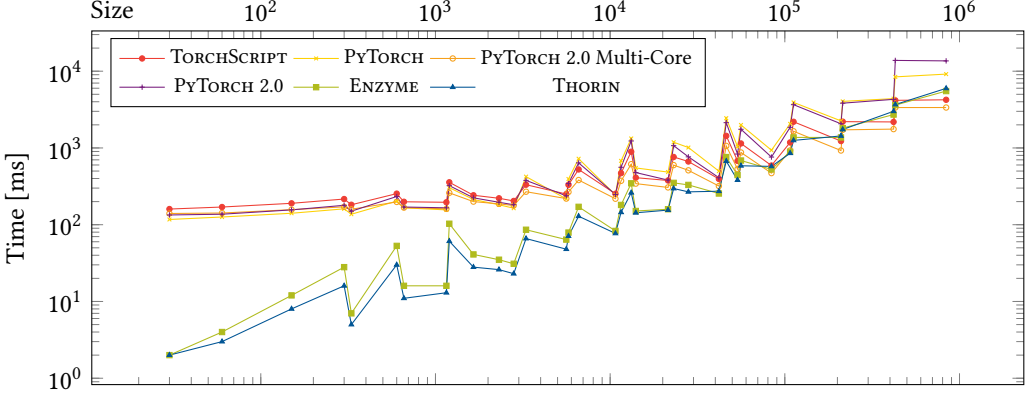


Fig. 8. Runtime of the Gaussian Mixture Model (GMM) algorithm on different input sizes (lower is better)

implementation<sup>5</sup> is written in Python and uses the PYTORCH package.

Our evaluation shows that ENZYME and %autodiff are comparable in performance. In some cases like bounded analysis, ENZYME has a better caching/recomputation balance resulting in lower runtime. In other cases like the GMM or LSTM benchmark, %autodiff is faster due to more caching. The main performance difference is due to caching or recomputation of intermediate results. ENZYME manages to better detect induction variables and recomputes them in the backward pass instead of storing them. On the other hand, it sometimes stores additional unnecessary intermediate results in the forward pass.

In our evaluation, PYTORCH is slower than ENZYME and %autodiff. This is partly due to the overhead of constructing the backward graph at runtime, the overhead of the Python interpreter, more memory usage, and missing optimizations like inlining. PYTORCH is inefficient in the bounded analysis benchmark resulting in a significantly longer running time. This slowdown is caused by the deeply nested function calls and loops that contain conditionals in the benchmark source code.

TORCHSCRIPT improves upon PYTORCH by compiling parts of the code. The resulting speedup is especially visible for very large input sizes where TORCHSCRIPT becomes one of the best implementations as shown in Figure 8. PYTORCH 2.0 further refines the TORCHSCRIPT approach of ahead-of-time compilation using TORCHDYNAMO [3], TORCHINDUCTOR, and AOTAUTOGRAD. For small inputs, the remaining overhead of Python, not compiled functions, and communication overhead still causes PYTORCH

2.0 to be slower than ENZYME and %autodiff. One advantage of PYTORCH is its support for utilizing multiple cores. In Figure 8, the PYTORCH 2.0 Multi-Core variant shows how this notably helps PYTORCH’s performance on large problem sizes. ENZYME and %autodiff both do not use multiple cores yet.

		Cyclomatic complexity		
		LoC [k]	total	avg.
ENZYME	w/	58.9	13,793	10.9
	w/o	49.4	11,857	12.8
%autodiff	w/	4.7	704	1.6
	w/o	1.3	141	1.8

w/o: without extensions & plugins  
w/: with extensions & plugins

Table 5. LoC and cyclomatic complexity measure (lower is better). We compute the measure per function; total cycl. complexity refers to the sum over all files.

<sup>5</sup><https://github.com/microsoft/ADBench/tree/master/src/python/modules/PyTorch>

6.3.2 *Code Complexity*. One major contribution of our approach to AD is its simplicity. In order to verify this claim, we estimate the code complexity of `%autodiff` and ENZYME using code complexity metrics. Table 5 summarizes the LoC and the *cyclomatic complexity* [36]. For ENZYME<sup>6</sup>, we measured the source folder. For better comparability, we also measured the metrics for ENZYME without the Clang plugin, the MLIR code, and without the different scalar evolution expander versions. The numbers without extensions refer to the plugin without the special casing of pointer arguments. ENZYME has roughly 10× more code than `%autodiff`. This does not necessarily give an indication of which code is simpler. But as a rule of thumb the larger a code base gets, the harder it becomes to debug and maintain. As a more profound code complexity metric, we also look at the *cyclomatic complexity* metrics<sup>7</sup> of both implementations. ENZYME’s cyclomatic complexity is roughly an order of magnitude larger than `%autodiff`’s.

## 6.4 Discussion on Dependent Type Checking

Programming with dependent types is often considered difficult as the programmer usually has to specify complex proof terms to persuade the type checker to accept their program. This essentially boils down to the type checker having to recognize that two expressions like  $f\ n$  and  $f\ m$  are the same. Now, with THORIN either normalization and partial evaluation are able to simplify both expressions to—let us say— $f\ o$ , or—if everything else fails—the programmer can simply insert a (potentially unsafe) cast via an unsafe axiom like `%core.bitcast` (see also Section 3.4).

## 7 RELATED & FUTURE WORK

THORIN lies at the intersection of higher-order, dependent type theory and low-level compiler IRs and is to the best of our knowledge the first of its kind in this regard. In this section, we will discuss work that influenced THORIN and related approaches.

*Partial Evaluation, AnyDSL & THORIN 1*. Partial Evaluation filters were first introduced by SCHISM [14] and are also used by the modern partial evaluation framework AnyDSL [32]. AnyDSL has been successfully applied for high-performance applications such as sequence alignment [41, 42] or ray tracing [46]. AnyDSL’s IR is THORIN 1 [33] while the THORIN version presented in this paper is actually THORIN 2 and a complete overhaul. THORIN 1 is based on CPS and has a set of built-in operations in direct style whereas THORIN 2 supports CPS and direct-style equally (see below). THORIN 1 lacks any of the advanced typing features of THORIN 2 such as polymorphism or dependent types and is not extensible via plugins. For this reason, AnyDSL relies on *shallow* DSL embedding and partial evaluation to remove any overhead. Like AnyDSL, THORIN 2 supports shallow DSL embedding but adds the possibility of *deep embeddings* via plugins (Section 2). Furthermore, THORIN 2 is to the best of our knowledge the first system that employs filter-based partial evaluation to resolve  $\beta$ -equivalence during type-checking a dependently typed language. While the idea of a “sea of nodes” goes back to Click and Paleczny [12], THORIN 1 pioneered this concept for higher-order languages. THORIN 2 inherits this representation and simplifies the graph even further: Due to THORIN 2’s roots in PTS, types—which are also just expressions—are part of the normal program graph as well.

*$\lambda$ -Calculi*. As a  $\lambda$ -calculus, THORIN takes inspiration from and shares similarities with many other  $\lambda$ -calculi—most notably: the  $\lambda$ -cube [7], PTS [37], CC [16], and the zip calculus [61]. Since THORIN is based upon a predicative flavor of CC it also subsumes a predicative flavor of System F [22, 49] and

<sup>6</sup><https://github.com/EnzymeAD/Enzyme>

<sup>7</sup>We used metrix++ for the measurement; see <https://github.com/metrixplusplus/metrixplusplus>

System  $F_\omega$ . Rossberg et al. [50] have shown that ML’s module system can be encoded in System  $F_\omega$ . Apart from impredicative idioms, THORIN can thus encode ML modules as well.

*CPS vs. Direct Style.* In the community for compilers of functional languages, there is a decades-old debate about whether to use CPS [4, 27, 29] or direct style [19, 34]. On the one hand, many optimizations such as simple rewrites like  $x + 0 \triangleright x$  are much easier to implement in direct-style. On the other hand, we need CPS to model unstructured control flow and at the end of the day, a compiled program consists of basic blocks with machine instructions that jump to each other—which is CPS. Thus, we second the opinion of Cong et al. [13] that the question of whether to compile with or without continuations is more a matter of what a compiler engineer/language designer wants to achieve and where in the compilation pipeline we are. THORIN itself does not really care whether you want to use CPS or direct style and provides syntactic sugar for both styles. We have also implemented a plugin `%direct` that lets you invoke direct-style functions as continuations and certain continuations as direct-style functions. For example, given `f` of type `.Fn T → U`, the expression `%direct.cps2ds f` has type `T → U`. Here we noted similarities to negation in the work of Ostermann et al. [43]. In addition, the `%direct` plugin CPS-converts direct-style function to CPS because THORIN’s LLVM backend expects CPS. Cong et al. [13] present a calculus that differentiates between first- and second-class continuations and a type system that ensures proper use. We have implemented another plugin for THORIN that does a similar classification via static analysis, instead. This plugin transforms escaping continuations via typed closure conversion [38] but is limited to non-dependent function types. Bowman and Ahmed [10] present a technique to closure-convert dependently typed functions in CC.

*Relationship to other IRs.* THORIN is flexible enough to mimic various IRs. We have already discussed in Section 3.2 how CPS makes THORIN akin to an SSA representation [5, 26]. THORIN can also model various extensions to SSA form such as (*thin*) *gated* SSA [23, 60], loop-closed SSA (LCSSA), or static single information (SSI) form [2]. It just depends on where additional variables are placed.

While MLIR [31] and THORIN are both extensible compiler frameworks and pursue similar goals, MLIR only provides a basic infrastructure for hosting languages that has “little builtin, everything customizable” [31, p. 3]. THORIN, on the other hand, aims to provide a general base language with an expressive type system. Similarly to THORIN’s plugins, MLIR offers extensibility through so-called *dialects*. Dialects enhance MLIR’s parser and type checker because MLIR lacks polymorphism or dependent types. This raises the question of how the different type-systems of the individual dialects interact and integrate. There are even dialects that violate basic SSA invariants (no cycles in the data dependence graph)<sup>8</sup>. Additionally, MLIR has only limited support for higher-order functions and relies on other dialects [8] to supply them. Finally, we argue that THORIN’s “sea of nodes”-style IR with only one syntactical element is easier to work with for high-level transformations than MLIR’s traditional compiler data structures (instruction lists, CFGs, ...).

Maziarz et al. [35] present an algorithm to hash expressions modulo  $\alpha$ -equivalence. This technique does not work for THORIN, however, as THORIN’s program graph is mutable at very specific spots (Section 4). Moreover, the *assignable* relation checks for compatible tuple types (Section 3.1.4) and resolves implicits (Section 5.2) anyway. By doing this, the relation additionally checks for  $\alpha$ -equivalence.

*Typed Assembly Language.* Previous work on low-level types [39] enriched imperative assembly languages with a type system. Building up on this work, Xi and Harper [64] added dependent types allowing for more advanced optimizations such as loop optimizations by utilizing array

<sup>8</sup><https://cirt.llvm.org/>

bounds. The underlying goal of utilizing additional information provided by the type system is the same as in our approach. However, we operate on a higher abstraction level with a functional intermediate representation instead of a low-level assembly level. As such, THORIN allows full dependent types with code generation via LLVM by eventual type-erasure. The extensible nature of THORIN furthermore makes the optimizations simpler allowing to represent all optimization stages in THORIN itself instead of having multiple separate compilation stages.

*Dependent Types.* Dependent types give more freedom in program construction and allow to express semantic properties in the types of expressions. Although they have been mainly used in proof assistants such as Coq, Agda, or Lean in the past, they have received some attention beyond proof assistants in recent years: Dependent Haskell [62] is an extension of Haskell that allows combining types and expressions while preserving backward compatibility with Haskell. Scala 3 is based upon dependent object types (DOT) [1] which features *path-dependent types*. These are a specific kind of dependent type where the dependent-upon value is a path.

Idris focuses on type-driven development but, if desired, properties of program behavior can be formally stated and proven. THORIN similarly expresses semantics on the type-level using dependent-types but differs in the focus on optimization instead of verification. Whereas Idris is a high-level programming language, THORIN as an IR is much closer to the hardware. Idris 2 [11] introduces Quantitative Type Theory (QTT) [6]. This not only allows to specify protocols as types but also to tag which types should be erased to clearly state which parts of a program materializes at runtime. THORIN on the other hand, uses partial evaluation filters to resolve  $\beta$ -equivalence and make this distinction.

F\* [56] combines automated theorem proving and interactive proving. The automation also aids when using refinement types as utilized in length-annotated arrays. THORIN’s user experience is similar as the normalization approach handles dependent types without further complications. However, THORIN does not aim to be a theorem prover. Furthermore, THORIN uses general dependent types instead of refinement types. This approach allows for more freedom as is exemplified in the AD plugin that computes the result type at the type level. Furthermore, the developer can extend the normalization approach to handle more complex cases. One possible rule could optimize a double reversion `rev (rev xs)` directly to `xs`; a use case, many automatic systems have problems with.

*Future Work.* In the future, we want to make it possible to specify normalizations directly as rewrite rules in THORIN’s surface language. This further reduces boilerplate C++ code. Moreover, this opens the door for additional sanity checks and non-deterministic rule sets that we could explore with rewrite engines like equality saturation [57]. For example, map/reduce-style rewrite rules have already been used with equality saturation to produce high-performance code [28]. We are also working on support for singleton, union, and intersection types. Singleton and union types are useful in many settings. Intersection types allow us to combine trait-like  $\Sigma$ -types such as `Cmp`  $\cap$  `Num`. We also want to add linear types or full QTT in order to validate by the type checker whether values that are supposed to be used linearly are used correctly. For example, values of type `%mem.M`—which track side-effects—must be used linearly, but this is right now not enforced by the type system. Furthermore, we want to enable parallelization and accelerators, such as GPUs, to enhance higher-level DSLs with these mechanisms. Finally, we want to reimplement existing DSLs such as Lift [55] or RISE/Shine [54] in THORIN.

## 8 CONCLUSION

In this paper, we presented THORIN, an extensible, higher-order intermediate representation. THORIN’s extensibility allows for expressing and optimizing programs at any level of abstraction. THORIN’s foundation in the the Calculus of Constructions provides a general and common type system for domain-specific languages to nest in. DSL authors benefit from reusing THORIN’s type system, normalization, and optimization framework without having to provide manual type-checkers and reimplementing standard optimizations for their custom operations. We have shown that using THORIN, we can generate code with state-of-the-art performance for high-level, domain-specific applications such as a RegEx matcher, automatic differentiation, as well as for low-level, imperative code.

## REFERENCES

- [1] Nada Amin, Samuel Grütter, Martin Odersky, Tiark Rompf, and Sandro Stucki. 2016. The Essence of Dependent Object Types. In *A List of Successes That Can Change the World - Essays Dedicated to Philip Wadler on the Occasion of His 60th Birthday (Lecture Notes in Computer Science, Vol. 9600)*, Sam Lindley, Conor McBride, Philip W. Trinder, and Donald Sannella (Eds.). Springer, 249–272. [https://doi.org/10.1007/978-3-319-30936-1\\_14](https://doi.org/10.1007/978-3-319-30936-1_14)
- [2] C. Scott Ananian. 1999. *The Static Single Information Form*. Ph. D. Dissertation. Princeton University.
- [3] Jason Ansel. 2022. *TorchDynamo*. <https://github.com/pytorch/torchdynamo>
- [4] Andrew W. Appel. 1992. *Compiling with Continuations*. Cambridge University Press.
- [5] Andrew W. Appel. 1998. SSA is Functional Programming. *ACM SIGPLAN Notices* 33, 4 (1998), 17–20. <https://doi.org/10.1145/278283.278285>
- [6] Robert Atkey. 2018. Syntax and Semantics of Quantitative Type Theory. In *Proceedings of the 33rd Annual ACM/IEEE Symposium on Logic in Computer Science, LICS 2018, Oxford, UK, July 09-12, 2018*, Anuj Dawar and Erich Grädel (Eds.). ACM, 56–65. <https://doi.org/10.1145/3209108.3209189>
- [7] Henk Barendregt. 1991. Introduction to Generalized Type Systems. *J. Funct. Program.* 1, 2 (1991), 125–154. <https://doi.org/10.1017/s0956796800020025>
- [8] Siddharth Bhat and Tobias Grosser. 2022. Lambda the Ultimate SSA: Optimizing Functional Programs in SSA. In *IEEE/ACM International Symposium on Code Generation and Optimization, CGO 2022, Seoul, Korea, Republic of, April 2-6, 2022*, Jae W. Lee, Sebastian Hack, and Tatiana Shpeisman (Eds.). IEEE, 1–11. <https://doi.org/10.1109/CGO53902.2022.9741279>
- [9] L Susan Blackford, Antoine Petitet, Roldan Pozo, Karin Remington, R Clint Whaley, James Demmel, Jack Dongarra, Iain Duff, Sven Hammarling, Greg Henry, et al. 2002. An updated set of basic linear algebra subprograms (BLAS). *ACM Trans. Math. Software* 28, 2 (2002), 135–151.
- [10] William J. Bowman and Amal Ahmed. 2018. Typed closure conversion for the calculus of constructions. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2018, Philadelphia, PA, USA, June 18-22, 2018*, 797–811. <https://doi.org/10.1145/3192366.3192372>
- [11] Edwin C. Brady. 2021. Idris 2: Quantitative Type Theory in Practice. In *35th European Conference on Object-Oriented Programming, ECOOP 2021, July 11-17, 2021, Aarhus, Denmark (Virtual Conference) (LIPIcs, Vol. 194)*, Anders Möller and Manu Sridharan (Eds.). Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 9:1–9:26. <https://doi.org/10.4230/LIPICS.ECOOP.2021.9>
- [12] Cliff Click and Michael Paleczny. 1995. A Simple Graph-Based Intermediate Representation. In *Proceedings ACM SIGPLAN Workshop on Intermediate Representations (IR’95), San Francisco, CA, USA, January 22, 1995*, 35–49. <https://doi.org/10.1145/202529.202534>
- [13] Youyou Cong, Leo Osvald, Grégory M. Essertel, and Tiark Rompf. 2019. Compiling with continuations, or without? whatever. *Proc. ACM Program. Lang.* 3, ICFP (2019), 79:1–79:28. <https://doi.org/10.1145/3341643>
- [14] Charles Consel. 1988. New Insights into Partial Evaluation: the SCHISM Experiment. In *ESOP ’88, 2nd European Symposium on Programming, Nancy, France, March 21-24, 1988, Proceedings (Lecture Notes in Computer Science, Vol. 300)*, Harald Ganzinger (Ed.). Springer, 236–246. [https://doi.org/10.1007/3-540-19027-9\\_16](https://doi.org/10.1007/3-540-19027-9_16)
- [15] Thierry Coquand and Jean Gallier. 1990. A Proof of Strong Normalization For the Theory of Constructions Using a Kripke-Like Interpretation. (07 1990).
- [16] Thierry Coquand and Gérard P. Huet. 1988. The Calculus of Constructions. *Inf. Comput.* 76, 2/3 (1988), 95–120. [https://doi.org/10.1016/0890-5401\(88\)90005-3](https://doi.org/10.1016/0890-5401(88)90005-3)
- [17] Li Deng. 2012. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine* 29, 6 (2012), 141–142.



- [18] Hana Dusíková. 2023. *Compile Time Regular Expression in C++*. <https://github.com/hanickadot/compile-time-regular-expressions>.
- [19] Cormac Flanagan, Amr Sabry, Bruce F. Duba, and Matthias Felleisen. 1993. The Essence of Compiling with Continuations. In *Proceedings of the ACM SIGPLAN'93 Conference on Programming Language Design and Implementation (PLDI), Albuquerque, New Mexico, USA, June 23-25, 1993*, Robert Cartwright (Ed.). ACM, 237–247. <https://doi.org/10.1145/155090.155113>
- [20] Free Software Foundation. 2023. `std::regex from libstdc++ version 13.1.1`. <https://gcc.gnu.org/onlinedocs/gcc-13.1.0/libstdc++/api/a07327.html>.
- [21] Syoyo Fujita. 2014. *aobench*. <https://code.google.com/archive/p/aobench/>.
- [22] Jean-Yves Girard. 1971. Une Extension De L'Interpretation De Gödel a L'Analyse, Et Son Application a L'Elimination Des Coupures Dans L'Analyse Et La Theorie Des Types. In *Proceedings of the Second Scandinavian Logic Symposium*, J.E. Fenstad (Ed.). Studies in Logic and the Foundations of Mathematics, Vol. 63. Elsevier, 63–92. [https://doi.org/10.1016/S0049-237X\(08\)70843-7](https://doi.org/10.1016/S0049-237X(08)70843-7)
- [23] Paul Havlak. 1993. Construction of Thinned Gated Single-Assignment Form. In *Languages and Compilers for Parallel Computing, 6th International Workshop, Portland, Oregon, USA, August 12-14, 1993, Proceedings (Lecture Notes in Computer Science, Vol. 768)*, Utpal Banerjee, David Gelernter, Alexandru Nicolau, and David A. Padua (Eds.). Springer, 477–499. [https://doi.org/10.1007/3-540-57659-2\\_28](https://doi.org/10.1007/3-540-57659-2_28)
- [24] Philip Hazel. 2021. *Perl-compatible Regular Expressions v2*. <http://www.pcre.org/current/doc/html/pcre2.html>.
- [25] John Hopcroft. 1971. An  $n \log n$  algorithm for minimizing states in a finite automaton. In *Theory of machines and computations*. Elsevier, 189–196.
- [26] Richard Kelsey. 1995. A Correspondence between Continuation Passing Style and Static Single Assignment Form. In *Proceedings ACM SIGPLAN Workshop on Intermediate Representations (IR'95), San Francisco, CA, USA, January 22, 1995*, Michael D. Ernst (Ed.). ACM, 13–23. <https://doi.org/10.1145/202529.202532>
- [27] Andrew Kennedy. 2007. Compiling with continuations, continued. In *Proceedings of the 12th ACM SIGPLAN International Conference on Functional Programming, ICFP 2007, Freiburg, Germany, October 1-3, 2007*, Ralf Hinze and Norman Ramsey (Eds.). ACM, 177–190. <https://doi.org/10.1145/1291151.1291179>
- [28] Thomas Koehler, Phil Trinder, and Michel Steuwer. 2021. Sketch-Guided Equality Saturation: Scaling Equality Saturation to Complex Optimizations in Languages with Bindings. *CoRR* abs/2111.13040 (2021). arXiv:2111.13040 <https://arxiv.org/abs/2111.13040>
- [29] David A. Kranz, Richard Kelsey, Jonathan Rees, Paul Hudak, and James Philbin. 1986. ORBIT: an optimizing compiler for scheme. In *Proceedings of the 1986 SIGPLAN Symposium on Compiler Construction, Palo Alto, California, USA, June 25-27, 1986*, Richard L. Wexelblat (Ed.). ACM, 219–233. <https://doi.org/10.1145/12276.13333>
- [30] Chris Lattner and Vikram S. Adve. 2004. LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation. In *2nd IEEE / ACM International Symposium on Code Generation and Optimization (CGO 2004), 20-24 March 2004, San Jose, CA, USA*. IEEE Computer Society, 75–88. <https://doi.org/10.1109/CGO.2004.1281665>
- [31] Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques A. Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, and Oleksandr Zinenko. 2021. MLIR: Scaling Compiler Infrastructure for Domain Specific Computation. In *IEEE/ACM International Symposium on Code Generation and Optimization, CGO 2021, Seoul, South Korea, February 27 - March 3, 2021*, Jae W. Lee, Mary Lou Soffa, and Ayal Zaks (Eds.). IEEE, 2–14. <https://doi.org/10.1109/CGO51591.2021.9370308>
- [32] Roland Leißa, Klaas Boesche, Sebastian Hack, Arsène Pérard-Gayot, Richard Membarth, Philipp Slusallek, André Müller, and Bertil Schmidt. 2018. AnyDSL: a partial evaluation framework for programming high-performance libraries. *Proc. ACM Program. Lang.* 2, OOPSLA (2018), 119:1–119:30. <https://doi.org/10.1145/3276489>
- [33] Roland Leißa, Marcel Köster, and Sebastian Hack. 2015. A graph-based higher-order intermediate representation. In *Proceedings of the 13th Annual IEEE/ACM International Symposium on Code Generation and Optimization, CGO 2015, San Francisco, CA, USA, February 07 - 11, 2015*. 202–212. <https://doi.org/10.1109/CGO.2015.7054200>
- [34] Luke Maurer, Paul Downen, Zena M. Ariola, and Simon L. Peyton Jones. 2017. Compiling without continuations. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2017, Barcelona, Spain, June 18-23, 2017*, Albert Cohen and Martin T. Vechev (Eds.). ACM, 482–494. <https://doi.org/10.1145/3062341.3062380>
- [35] Krzysztof Maziarz, Tom Ellis, Alan Lawrence, Andrew W. Fitzgibbon, and Simon Peyton Jones. 2021. Hashing modulo alpha-equivalence. In *PLDI '21: 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation, Virtual Event, Canada, June 20-25, 2021*, Stephen N. Freund and Eran Yahav (Eds.). ACM, 960–973. <https://doi.org/10.1145/3453483.3454088>
- [36] Thomas J McCabe. 1976. A complexity measure. *IEEE Transactions on software Engineering* 4 (1976), 308–320.
- [37] E Meijer and Simon Peyton Jones. 1997. *Henk: a typed intermediate language* (types in compilation ed.). <https://www.microsoft.com/en-us/research/publication/henk-a-typed-intermediate-language/>



- [38] Yasuhiko Minamide, J. Gregory Morrisett, and Robert Harper. 1996. Typed Closure Conversion. In *Conference Record of POPL '96: The 23rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Papers Presented at the Symposium, St. Petersburg Beach, Florida, USA, January 21-24, 1996*. 271–283. <https://doi.org/10.1145/237721.237791>
- [39] J. Gregory Morrisett, David Walker, Karl Crary, and Neal Glew. 1998. From System F to Typed Assembly Language. In *POPL '98, Proceedings of the 25th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, San Diego, CA, USA, January 19-21, 1998*, David B. MacQueen and Luca Cardelli (Eds.). ACM, 85–97. <https://doi.org/10.1145/268946.268954>
- [40] William Moses and Valentin Churavy. 2020. Instead of rewriting foreign code for machine learning, automatically synthesize fast gradients. *Advances in neural information processing systems* 33 (2020), 12472–12485.
- [41] André Müller, Bertil Schmidt, Andreas Hildebrandt, Richard Membarth, Roland Leißa, Matthias Kruse, and Sebastian Hack. 2020. AnySeq: A High Performance Sequence Alignment Library based on Partial Evaluation. In *2020 IEEE International Parallel and Distributed Processing Symposium (IPDPS), New Orleans, LA, USA, May 18-22, 2020*. IEEE, 1030–1040. <https://doi.org/10.1109/IPDPS47924.2020.00109>
- [42] André Müller, Bertil Schmidt, Richard Membarth, Roland Leißa, and Sebastian Hack. 2022. AnySeq/GPU: a novel approach for faster sequence alignment on GPUs. In *ICS '22: 2022 International Conference on Supercomputing, Virtual Event, June 28 - 30, 2022*, Lawrence Rauchwerger, Kirk W. Cameron, Dimitrios S. Nikolopoulos, and Dionisios N. Pnevmatikatos (Eds.). ACM, 20:1–20:11. <https://doi.org/10.1145/3524059.3532376>
- [43] Klaus Ostermann, David Binder, Ingo Skupin, Tim Süßerkrüb, and Paul Downen. 2022. Introduction and elimination, left and right. *Proc. ACM Program. Lang.* 6, ICFP (2022), 438–465. <https://doi.org/10.1145/3547637>
- [44] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.). Curran Associates, Inc., 8024–8035. <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
- [45] Barak A Pearlmutter and Jeffrey Mark Siskind. 2008. Reverse-mode AD in a functional framework: Lambda the ultimate backpropagator. *ACM Transactions on Programming Languages and Systems (TOPLAS)* 30, 2 (2008), 1–36.
- [46] Arsène Pérard-Gayot, Richard Membarth, Roland Leißa, Sebastian Hack, and Philipp Slusallek. 2019. Rodent: generating renderers without writing a generator. *ACM Trans. Graph.* 38, 4 (2019), 40:1–40:12. <https://doi.org/10.1145/3306346.3322955>
- [47] Randy Pollack and Eric Poll. 1992. Typechecking in Pure Type Systems. In *Informal proceedings of Logical Frameworks'92*. 271–288.
- [48] Dragomir Radev. 2008. *CLAIR collection of fraud email*.
- [49] John C. Reynolds. 1974. Towards a theory of type structure. In *Programming Symposium, Proceedings Colloque sur la Programmation, Paris, France, April 9-11, 1974 (Lecture Notes in Computer Science, Vol. 19)*, Bernard J. Robinet (Ed.). Springer, 408–423. [https://doi.org/10.1007/3-540-06859-7\\_148](https://doi.org/10.1007/3-540-06859-7_148)
- [50] Andreas Rossberg, Claudio V. Russo, and Derek Dreyer. 2014. F-ing modules. *J. Funct. Program.* 24, 5 (2014), 529–607. <https://doi.org/10.1017/S0956796814000264>
- [51] Bhargav Shivkumar, Jeffrey C. Murphy, and Lukasz Ziarek. 2021. Real-time MLton: A Standard ML runtime for real-time functional programs. *J. Funct. Program.* 31 (2021), e19. <https://doi.org/10.1017/S0956796821000174>
- [52] Matthieu Sozeau, Simon Boulier, Yannick Forster, Nicolas Tabareau, and Théo Winterhalter. 2020. Coq Coq correct! verification of type checking and erasure for Coq, in Coq. *Proc. ACM Program. Lang.* 4, POPL (2020), 8:1–8:28. <https://doi.org/10.1145/3371076>
- [53] Filip Srajer, Zuzana Kukelova, and Andrew Fitzgibbon. 2018. A benchmark of selected algorithmic differentiation tools on some problems in computer vision and machine learning. *Optimization Methods and Software* 33, 4-6 (2018), 889–906.
- [54] Michel Steuwer, Thomas Koehler, Bastian Köpcke, and Federico Pizzuti. 2022. RISE & Shine: Language-Oriented Compiler Design. *CoRR abs/2201.03611* (2022). arXiv:2201.03611 <https://arxiv.org/abs/2201.03611>
- [55] Michel Steuwer, Toomas Rimmelg, and Christophe Dubach. 2017. Lift: a functional data-parallel IR for high-performance GPU code generation. In *Proceedings of the 2017 International Symposium on Code Generation and Optimization, CGO 2017, Austin, TX, USA, February 4-8, 2017*, Vijay Janapa Reddi, Aaron Smith, and Lingjia Tang (Eds.). ACM, 74–85. <http://dl.acm.org/citation.cfm?id=3049841>
- [56] Nikhil Swamy, Catalin Hritcu, Chantal Keller, Aseem Rastogi, Antoine Delignat-Lavaud, Simon Forest, Karthikeyan Bhargavan, Cédric Fournet, Pierre-Yves Strub, Markulf Kohlweiss, Jean Karim Zinzindohoue, and Santiago Zanella Béguelin. 2016. Dependent types and multi-monadic effects in F. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2016, St. Petersburg, FL, USA, January 20 - 22, 2016*, Rastislav Bodik and Rupak Majumdar (Eds.). ACM, 256–270. <https://doi.org/10.1145/2837614.2837655>

- [57] Ross Tate, Michael Stepp, Zachary Tatlock, and Sorin Lerner. 2009. Equality saturation: a new approach to optimization. In *Proceedings of the 36th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2009, Savannah, GA, USA, January 21-23, 2009*, Zhong Shao and Benjamin C. Pierce (Eds.). ACM, 264–276. <https://doi.org/10.1145/1480881.1480915>
- [58] Benchmarksgame Team. 2023. *The Computer Language Benchmarks Game*. <https://benchmarksgame-team.pages.debian.net/benchmarksgame/>.
- [59] Andrew P. Tolmach and Dino Oliva. 1998. From ML to Ada: Strongly-Typed Language Interoperability via Source Translation. *J. Funct. Program.* 8, 4 (1998), 367–412. <https://doi.org/10.1017/s0956796898003086>
- [60] Peng Tu and David A. Padua. 1995. Efficient Building and Placing of Gating Functions. In *Proceedings of the ACM SIGPLAN'95 Conference on Programming Language Design and Implementation (PLDI), La Jolla, California, USA, June 18-21, 1995*, David W. Wall (Ed.). ACM, 47–55. <https://doi.org/10.1145/207110.207115>
- [61] Mark Tullsen. 2000. The Zip Calculus. In *Mathematics of Program Construction, 5th International Conference, MPC 2000, Ponte de Lima, Portugal, July 3-5, 2000, Proceedings*. 28–44. [https://doi.org/10.1007/10722010\\_3](https://doi.org/10.1007/10722010_3)
- [62] Stephanie Weirich, Antoine Voizard, Pedro Henrique Azevedo de Amorim, and Richard A. Eisenberg. 2017. A specification for dependent types in Haskell. *Proc. ACM Program. Lang.* 1, ICFP (2017), 31:1–31:29. <https://doi.org/10.1145/3110275>
- [63] Reinhard Wilhelm, Helmut Seidl, and Sebastian Hack. 2013. *Compiler Design - Syntactic and Semantic Analysis*. Springer. <https://doi.org/10.1007/978-3-642-17540-4>
- [64] Hongwei Xi and Robert Harper. 2001. A Dependently Typed Assembly Language. In *Proceedings of the Sixth ACM SIGPLAN International Conference on Functional Programming (ICFP '01), Firenze (Florence), Italy, September 3-5, 2001*, Benjamin C. Pierce (Ed.). ACM, 169–180. <https://doi.org/10.1145/507635.507657>