

# ScalaCT: Type-Directed Staging at Compile-Time

Vojin Jovanovic   Eugene Burmako   Denys Shabalin

École polytechnique fédérale de Lausanne—EPFL

{first}.{last}@epfl.ch

## Abstract

Staging systems choose a compilation stage in which a term is executed based on: quotation (e.g., MetaOCaml), or types (e.g., LMS). Type based staging systems, require type annotations of all future stage terms, as well as implementing reification and code generation logic for all future stage types. Further, when we use staging at host language compile-time, all terms scheduled to execute at runtime require type annotations and all libraries used at runtime require logic for reification and code generation. Number of annotations is especially noticeable in languages with local type inference as method parameter types must be provided by the programmer.

We introduce a type based staging system for Scala where terms whose types are annotated are executed in the earlier stage of compilation, in our case, at host language compile-time. Annotated types represent merely compile-time views of original types and therefore no reification and code generation logic is necessary. We compare our framework with LMS and show that in majority of programs we require less type annotations while we achieve same performance improvements.

**Categories and Subject Descriptors** D.3.3 [Programming Languages]: Language Constructs and Features

**Keywords** Multi-Stage Programming, Partial Evaluation

## 1. Introduction

*Multi-stage programming* (or *staging*) is a meta-programming technique where compilation is separated in multiple *stages*. Execution of each stage outputs code that is executed in the next stage of compilation. The first stage of compilation happens at the *host language* compile time, the second

stage happens at the host language runtime, the third stage happens at runtime of runtime generated code, etc. Different stages of compilation can be executed in the same language [22, 34] or in different languages [4, 11]. In this work we will focus on staging systems where all stages are in the same language and that, through static typing, assure that terms in the next stage are well typed.

Notable staging systems in statically typed languages are MetaOCaml [5, 34] and LMS [25]. These systems were successfully applied as a *partial evaluator* [16]: for removing abstraction overheads in high-level programs [6, 25], for domain-specific languages [7, 17, 33], and for converting language interpreters into compilers [12, 27]. Staging originates from research on two-level [9, 22] and multi-level [10] calculi.

We show an example of how staging is used for partial evaluation of a function for computing the inner product of two vectors<sup>1</sup>:

```
def dot[T:Numeric](v1: Vector[T], v2: Vector[T]): T =  
  (v1 zip v2).foldLeft(zero[T]) {  
    case (prod, (c1, cr)) => prod + c1 * cr  
  }
```

In function `dot`, if vector sizes are constant, the inner product can be partially evaluated into a sum of products of vector components. To achieve partial evaluation, we must communicate to the staging system that operations on values of vector components should be executed in the next stage. The compilation stage in which a term is executed is determined by *code quotation* (in MetaOCaml) or by parametric types `Rep` (in LMS). In LMS marking that the vector size is statically known is achieved by annotating only vector elements with a `Rep` type<sup>2</sup>:

```
def dot[T:Numeric]  
  (v1: Vector[Rep[T]], v2: Vector[Rep[T]]): Rep[T]
```

Here the `Rep` annotations on `Rep[T]` denote that elements of vectors will be known only in the next stage (in LMS, this is a stage after run-time compilation). After run-time compilation `zip`, `foldLeft`, and pattern matching inside the closure will not exist as they were evaluated in the previous stage of compilation (host language runtime). Note that in LMS unannotated code is always executed during

[Copyright notice will appear here once 'preprint' option is removed.]

host-language runtime and type-annotated code is executed after run-time compilation.

**Staging at host language compile time.** How can we use staging for programs whose values are statically known at the host language compile-time (the first stage)? Existing staging frameworks treat unannotated terms as runtime values of the host language and annotated terms as values in later stages of compilation. Even if we would take that the first stage is executed at the host language compile time, we would have to annotate all run-time values. Annotating all values is cumbersome since host language run-time values comprise the majority of user programs (§4).

Interesting exception is MacroML

**Code Duplication.** Staging systems based on type annotations (e.g., LMS and type-directed partial evaluation [8]) inherently require code duplication as, a priori, no operations are defined on `Rep` annotated types. For example, in the LMS version of the `dot` function, all numeric types (i.e., `Rep[Int]`, `Rep[Double]`, etc.) must be re-implemented in order to typecheck the programs and achieve code generation for the next stage.

Sujeeth et al. [31] and Jovanovic et al. [18] propose generating code for the next stage computations based on a language specification. These approaches solve the problem, but they require writing additional specification for the libraries, require a large machinery for code generation, and support only restricted parts of Scala.

**Annotating the Previous Stage.** The main idea of this paper is that *annotated types* should denote computations that happen during the *previous stage* of compilation. The reason is that static terms appear less frequently than run-time terms in a large set of analyzed programs (§4). Therefore, annotating static terms introduces less overhead for the programmer.

We treat annotated types as *compile-time views* of existing data types. Compile-time view of a type denotes that all operations on that type are executed at host language compile time. We promote types to their compile-time views with the `@ct` annotation (e.g., `Int@ct`). Similarly, statically known terms can be promoted their compile time duals with the `ct` function on the term level. By having two views of the same type we obviate the need for introducing reification and code generation logic for existing types.

With compile-time views, to require that vectors `v1` and `v2` are static and to partially evaluate the function, a programmer needs to make a simple modification of the `dot` signature:

```
def dot[V: Numeric@ct]
  (v1: Vector[V]@ct, v2: Vector[V]@ct): V
```

<sup>1</sup> All code examples are written in *Scala*. It is necessary to know the basics of Scala to comprehend this paper.

<sup>2</sup> In this work we use LMS as a representative of type-based staging systems.

Since, vector elements are polymorphic the result of the function can be a dynamic value, or a compile-time view that can be further used for compile-time computations. The binding time of the return type of `dot` will match the binding time of vector elements:

```
// [e11, e12, e13, e14] are dynamic decimals
dot(Vector(e11, e12), Vector(e13, e14))
  ⇨ (e11 * e13 + e12 * e14): Double

// static terms are internally tracked
dot(Vector(2.0, 4.0), Vector(1.0, 10.0))
  ⇨ (2.0 * 1.0 + 4.0 * 10.0): Double@static

// ct promotes static terms to compile-time views
dot(Vector(ct(2), ct(4)),
  Vector(ct(1), ct(10)))
  ⇨ 42: Double@ct
```

In this paper we contribute to the state of the art:

- By introducing compile-time views (§2) as means to succinctly achieve type safe two-stage programming starting from host language compile time.
- By obviating the need for reification and code generation logic in type based staging systems.
- By demonstrating the usefulness of compile-time views in four case studies (§3): inlining, partially evaluating recursion, removing overheads of variable argument functions, and removing overheads of type-classes [15, 24, 35].

We have implemented a staging extension for Scala `ScalaCT`<sup>3</sup>. `ScalaCT` has a minimal interface (§2) based on type annotations. We have evaluated performance gains and the validity of `ScalaCT` on all case studies (§3) and compared them to LMS. In all benchmarks (§5) our evaluator performs the same as LMS and gives significant performance gains compared to original programs.

## 2. Compile-Time Views in Scala

In this section we informally present `ScalaCT`, a staging extension for Scala based on compile-time views. `ScalaCT` is a compiler plugin that executes in a phase after the Scala type checker. The plugin takes as input typechecked Scala programs and uses type annotations [23] to track and verify information about the binding-time of terms. It supports only two stages of compilation: host language compile-time (types annotated with `@ct`) and host language run-time (unannotated code).

To the user, `ScalaCT` exposes a minimal interface (Figure 2) with a single annotation `ct` and a single function `ct`.

**Annotation** `ct` is used on types (e.g., `Int@ct`) to promote them to their compile-time views. The annotation is integrated in the Scala’s type system and, therefore, can be arbitrarily nested in different variants of types.

<sup>3</sup>Source code: <https://github.com/scala-inline/scala-inline>.

```
package object scalact {
  final class ct extends StaticAnnotation

  @compileTimeOnly def ct[T](body: => T): T = ???
}
```

**Figure 1.** Interface of ScalaCT.

Since all operations on compile-time views are executed at compile time, non-generic method parameters and result types of compile-time views also become compile-time views. Table 1 shows how the `@ct` annotation can be placed on types and how it affects method signatures on annotated types.

In Table 1, on `Int@ct` both parameters and result types of all methods are also compile-time views. On the other hand, `Vector[Int]@ct` has parameters of all methods transformed except the generic ones. In effect, this, makes higher order combinators of `Vector` operate on dynamic values, thus, function `f` passed to `map` accepts the dynamic value as input. Type `Vector[Int@ct]@ct` has all methods executed at compile-time. The return type of the function `map` on `Vector[Int@ct]@ct` can still be either dynamic or a compile-time view due to the type parameter `U`.

Annotation `ct` can be used to achieve simple inlining of statically known methods and functions. This is achieved by putting the annotation of the method/function definition:

```
def dot[V: Numeric]
  (v1: Vector[V], v2: Vector[V]): V
```

Annotated methods will have an annotated method type

```
((v1: Vector[V], v2: Vector[V]) => V)@ct
```

which can not be written by the users. This is not the first time that inlining is achieved through partial evaluation [20].

**Function** `ct` is used at the term level for promoting literals, modules, and methods/functions into their compile-time views. Without `ct` we would not be able to instantiate compile-time views of types. Table 2 shows how different types of terms are promoted to their compile-time views. An exception for promoting terms to compile-time views is the `new` construct. Here we use the type annotation on the constructed type.

## 2.1 Tracking Binding-Time of Terms

Internally ScalaCT has additional type annotations for tracking the binding time of terms. Type of each term is annotated with either `dynamic`, `static`, or `ct`. `dynamic` denotes that the term can only be known at runtime, `static` that the term is known at compile-time but it will not be computed at compile time, and `ct` that the term will be computed at compile-time.

Tracking static terms was studied in the context of binding-time analysis in partial evaluation for typed [21] and untyped [14] languages. We use similar techniques, however, unlike in partial evaluation we do not evaluate static terms at compile time. They are tracked for verify-

ing correctness and providing convenient implicit conversions. Static terms are evaluated only when they are explicitly marked by the programmer with `ct`.

In ScalaCT language literals, functions, direct class constructor calls with static arguments, and static method calls with static arguments are marked as static. Examples of static terms are

```
1.0, "1", (x: Int => x), new Cons(1, Nil), List(1,2,3)
```

## 2.2 Least Upper Bounds

We use subtyping of Scala to simplify tracking of binding times by introducing a subtyping relation between `dynamic`, `static`, and `ct`. We argue that a `static` type is a more specific `dynamic` as it is statically known and that `ct` is more specific than `static` as its operations are executed at compile time. Therefore we establish that

```
ct <: static <: dynamic
```

The use of subtyping simplifies tracking binding times of terms as in all cases where least upper bounds are calculated we can use the same mechanism for binding-times. An interesting example are the binding times of type parameters:

```
ct(List)(1, ct(2)): List[Int@static]@ct
ct(List)(ct(1), ct(2)): List[Int@ct]@ct
ct(List)((x: Int@dynamic), ct(2)): List[Int@dynamic]@ct
```

Notable exception are control flow constructs for which the original Scala rules for least upper bounds do not hold. The binding-time of control flow constructs does not depend only on the return type of the branches but also on the conditionals. For example, if both branches of an `if` construct are `static` the result can still be `dynamic` if the condition is `dynamic`. Here subtyping also helps as the binding type of the return value is simply calculated as `lub(c, thn, elz)` where `lub(tps: Type*)` is a function for computing the least upper bounds of types, and `c`, `thn`, `elz` are respectively binding times of the condition, the then branch, and the else branch. The same principles can be applied for pattern matching.

## 2.3 Well-Formedness of Compile-Time Views

Earlier stages of computation can not depend on values from later stages. This property, defined as *cross-stage persistence* [34, 36], imposes that all operations on compile-time views must known at compile time.

To satisfy cross-stage persistence ScalaCT verifies that binding time of composite types (e.g., polymorphic types, function types, record types, etc.) is always a subtype of the binding time of their components. In the following example, we show malformed types and examples of terms that are inconsistent:

```
xs: List[Int@ct]      => ct(Predef).println(xs.head)
fn: (Int@ct=>Int@ct) => ct(Predef).println(fn(ct(1)))
```

**Table 1.** Compile-time views of types and their corresponding method signatures.

Annotated Type	Type's Method Signatures
<code>Int@ct</code>	<code>+(rhs: Int@ct): Int@ct</code>
<code>Vector[Int]@ct</code>	<code>map[U](f: (Int =&gt; U)@ct): Vector[U]@ct</code> <code>length: Int@ct</code>
<code>Vector[Int@ct]@ct</code>	<code>map[U](f: (Int@ct =&gt; U)@ct): Vector[U]@ct</code> <code>length: Int@ct</code>
<code>Map[Int@ct, Int]@ct</code>	<code>get(key: Int@ct): Option[Int]@ct</code>

**Table 2.** Promotion of terms to their compile-time views.

Promoted Term	Term's Type
<code>ct(Vector)(1, 2, 3)</code>	<code>: Vector[Int]@ct</code>
<code>ct(Vector)(ct(1), ct(2), ct(3))</code>	<code>: Vector[Int@ct]@ct</code>
<code>ct((x: Int@ct) =&gt; x)</code>	<code>: (Int@ct =&gt; Int@ct)@ct</code>
<code>ct((x: Int) =&gt; x)</code>	<code>: (Int =&gt; Int)@ct</code>
<code>new (::@ct)(1, Nil)</code>	<code>: (::[Int])@ct</code>
<code>new (::@ct)(ct(1), ct(Nil))</code>	<code>: (::[Int@ct])@ct</code>

In the first example the program would, according to the semantics of `@ct`, print a head of the list at compile time. However, the head of the list is known only in the runtime stage. In the second example the program should print the result of `fn` at compile time but the body of the function will be known only at runtime. By causality such examples are not possible.

On functions/methods the `ct` annotation requires that function/method bodies are known at compile-time. Otherwise, inlining of such functions/methods would not be possible at compile-time. In Scala, method bodies are statically known in objects and classes with final methods, thus, the `ct` annotation is only applicable on such methods.

## 2.4 Implicit Conversions

If method parameters require compile-time views of a type the corresponding arguments in method application would always have to be promoted to `ct`. In some libraries this could require an inconveniently large number of annotations.

To minimize the number of required annotations we introduce implicit conversions from certain `static` terms to `ct` terms. The conversions support translation of language literals, direct class constructor calls with static arguments, and static method calls with static arguments into their compile-time views. Since our compile-time evaluator does not use Asai's [2, 32] method to keep track of the value of each static term, we disallow implicit conversions of terms with static variables.

For example, for a factorial function

```
def fact(n: Int @ct): Int@ct =
  if (n == 0) 1 else fact(n - 1)
```

we will not require promotions of literals 0, and 1. Furthermore, the function can be invoked without promoting the argument into its compile-time view:

```
fact(5)
  ↪ 120
```

Without implicit conversions the factorial functions would be more verbose

```
def fact(n: Int @ct): Int@ct =
  if (n == ct(0)) ct(1) else fact(n - ct(1))
```

as well as each function application (`fact(ct(5))`).

## 3. Case Studies

In this section we present selected use-cases for compile-time views that, at the same time, demonstrate step-by-step the mechanics behind ScalaCT. We start by inlining a simple function with staging (§3.1), then do the canonical staging example of the integer power function (§3.2), then we demonstrate how variable argument functions can be desugared into the core functionality (§3.3). Finally, we demonstrate how the abstraction overhead of the `dot` function and all associated type-class related abstraction can be removed (§3.5). For formal partial evaluation rules refer c.f. §??.

### 3.1 Inlining Expressed Through Staging

Function inlining can be expressed as staged computation [20]. Inlining is achieved when a statically known function body is applied with symbolic arguments. In ScalaCT we use the `ct` annotation on functions and methods to achieve inlining:

```
@ct def zero[T](implicit num: Numeric[T]) = num.zero

zero[Double]
  ↪ num.zero
```

### 3.2 Recursion

The canonical example in staging literature is partial evaluation of the power function where exponent is an integer:

```
def pow(base: Double, exp: Int): Double =
  if (exp == 0) 1 else base * pow(base, exp - 1)
```

When the exponent (`exp`) is statically known this function can be partially evaluated into `exp` multiplications of the base argument, significantly improving performance [5].

With compile-time views making `pow` partially evaluated requires adding only one annotation:

```
def pow(base: Double, exp: Int@ct): Double =
  if (exp == 0) 1 else base * pow(base, exp - 1)
```

To satisfy cross-stage persistence (§2.3) the `pow` must be `@ct`. However, to reduce the number of required annotations we implicitly add the `ct` annotation when at least one parameter type or the result type is marked as `ct`. In the example the `ct` annotation on `exp` requires that the function must be called with a compile-time view of `Int`. `ScalaCT` ensures that the definition of the `pow` function does not cause infinite recursion at compile-time by invoking the power function only when the value of the `ct` arguments is known.

The application of the function `pow` with a constant exponent will produce:

```
pow(base, 4)
  ↪ base * base * base * base * 1
```

Constant 4 is promoted to `ct` by the implicit conversions (§2.4).

### 3.3 Variable Argument Functions

Variable argument functions appear in widely used languages like Java, C#, and Scala. Such arguments are typically passed in the function body inside of the data structure (e.g. `Seq[T]` in Scala). When applied with variable arguments the size of the data-structure is statically known and all operations on them can be partially evaluated. However, sometimes, the function is called with arguments of dynamic size. For example, function `min` that accepts multiple integers

```
def min(vs: Int*): Int = vs.tail.foldLeft(vs.head) {
  (min, el) => if (el < min) el else min
}
```

can be called either with statically known arguments (e.g., `min(1,2)`) or with dynamic arguments:

```
val values: Seq[Int] = ... // dynamic value
min(values: _*)
```

Ideally, we would be able to achieve partial evaluation if the arguments are of statically known size and avoid partial evaluation in case of dynamic arguments. To this end we translate the method `min` into a partially evaluated version and a dynamic version. The call to these methods is dispatched, at compile-time, by the `min` method which checks

```
def min(vs: Int*): Int = macro
  if (isVarargs(vs)) q"min_CT(vs)"
  else q"min_D(vs)"

def min_CT(vs: Seq[Int]@ct): Int =
  vs.tail.foldLeft(vs.head) { (min, el) =>
    if (el < min) el else min
  }

def min_D(vs: Seq[Int]): Int =
  vs.tail.foldLeft(vs.head) {
    (min, el) => if (el < min) el else min
  }
```

**Figure 2.** Function `min` is desugared into a `min` macro that based on the binding time of the arguments dispatches to the partially evaluated version (`min_CT`) for statically known varargs or to the original `min` function for dynamic arguments `min_D`.

```
object Numeric {
  implicit def dnum: Numeric[Double]@ct =
    ct(DoubleNumeric)
  def zero[T](implicit num: Numeric[T]@ct): T =
    num.zero
}

trait Numeric[T] {
  def plus(x: T, y: T): T
  def times(x: T, y: T): T
  def zero: T
}

object DoubleNumeric extends Numeric[Double] {
  def plus(x: Double, y: Double): Double = x + y
  def times(x: Double, y: Double): Double = x * y
  def zero: Double = 0.0
}
```

**Figure 3.** Removing abstraction overheads of type classes.

if arguments are statically known. Desugaring of `min` is shown in Figure 2.

### 3.4 Removing Abstraction Overhead of Type-Classes

Type-classes are omnipresent in everyday programming as they allow abstraction over generic parameters (e.g., `Numeric` abstracts over numeric values). Unfortunately, type-classes introduce *dynamic dispatch* on every call [26] and, thus, impose a performance penalty. Type-classes are in most of the cases statically known. Here we show how with `ScalaCT` we can remove all abstraction overheads of type classes.

In Scala, type classes are implemented with objects and implicit parameters [24]. In Figure 3, we define a trait `Numeric` serves as an interface for all numeric types. Then we define a concrete implementation of `Numeric` for type `Double` (`DoubleNumeric`). The `DoubleNumeric` is then passed as an implicit argument `dnum` to all methods that use it (e.g., `zero`).

When `zero` is applied first the implicit argument (`dnum`) gets inlined due to the `ct` annotation of the return type,

then the function `zero` gets inlined. Since `dnum` returns a compile-time view of `DoubleNumeric` the method `zero` on `dnum` is evaluated at compile time. The constant `0.0` is promoted to `ct` since `DoubleNumeric` is a compile time view. Finally the `ct(0.0)` result is coerced to a dynamic value by the signature of `Numeric.zero`. The compile-time execution is shown in the following snippet

```
Numeric.zero[Double]
  ↪ Numeric.zero[Double](DoubleNumeric)
  ↪ ct(DoubleNumeric).zero
  ↪ (ct(0.0): Double)
  ↪ 0.0
```

### 3.5 Inner Product of Vectors

Here we demonstrate how the introductory example (§1) is partially evaluated through staging. We start with the desugared `dot` function (i.e., all implicit operations are shown):

```
def dot[V](v1: Vector[V]@ct, v2: Vector[V]@ct)
  (implicit num: Numeric[V]@ct): V =
  (v1 zip v2).foldLeft(zero[V](num)) {
    case (prod, (c1, cr)) => prod + c1 * cr
  }
```

Function `dot` is generic in the type of vector elements. This will reflect upon the staging annotations as well (`ct` and `static`). When we apply the `dot` function with static arguments we will get the vector with static elements back:

```
dot[Double@static](
  ct(Vector)(2.0, 4.0), ct(Vector)(1.0, 10.0))(
  Numeric.dnum)
  ↪
  (ct(Vector)(2.0, 4.0) zip ct(Vector)(1.0, 10.0))
  .foldLeft(ct(0.0)) {
    case (prod, (c1, cr)) => prod + c1 * cr
  }
  ↪ (2.0 * 1.0 + 4.0 * 10.0): Double@static
```

When `dot` is evaluated with the `ct` elements the last step will further execute to a single compile-time value that can further be used in compile-time computations:

```
dot[Double@ct](
  ct(Vector)(ct(2.0), ct(4.0)),
  ct(Vector)(ct(1.0), ct(10.0)))(Numeric.dnum)
  ↪ ct(2.0) * ct(1.0) + ct(4.0) * ct(10.0)
  ↪ 42.0: Double@ct
```

## 4. Discussion

To distinguish terms executed at compile-time from terms executed at runtime with type annotations we have the following possibilities:

1. Annotate types of all terms that should be executed at runtime. Here all types analyzed LMS and realized that this is not an option.
2. Annotate types of terms that should be executed at runtime but introduce scopes (e.g., method bodies) for which this rule applies. In this way we would avoid annotating types of all runtime terms. This approach is taken by

MacroML where macro functions are executed at compile time and quoted terms are executed at runtime. First approach is, also, a special case of this approach where there is a single scope for the whole language.

3. Annotate types of terms that are executed at compile time. This approach is used with ScalaCT and annotated types are called compile-time views.

In ScalaCT we decided to annotate types of terms that are executed at compile time. Compared to the first solution our approach takes requires less annotations. We analyzed

Compared to the second approach our solution is simpler to comprehend and communicate. In the second approach there are two things that users need to understand when reasoning about staged programs: *i)* where does the compile time scope start, and *ii)* which terms are annotated. With ScalaCT the comprehension is simple: terms whose types are annotated with `ct` are executed at compile time.

It is also interesting to the second and third approaches. Here the number of annotations depends on the program. If the programs are mostly partially evaluated the second approach is better. These category of programs could also be regarded as code generators as most of the code is executed at compile time and produces large outputs. When programs are comprised of mostly runtime values the approach of ScalaCT requires less annotations.

## 5. Evaluation

In this section we evaluate the amount of code that is obviated compared to existing type directed staging systems (§5.1). Then we evaluate performance of ScalaCT compared to LMS and hand optimized code (§5.2)

### 5.1 Reduction in Code Duplication

Evaluating reduction of duplicated code (for reification and code generation) in type based staging systems is difficult as the factor varies from program to program. To avoid benchmark dependent results we instead calculate the lower bound on the duplication factor.

Given that we have a method on a type `T` whose body contains `n` lines of code (without the method definition). To introduce the same method on an annotated type `Rep[T]` we need another method for reification which has at least 1 line of code. Then we need code generation logic, which, if we use the same language should not have less lines than the original method plus at least one line for matching the reified method. For method of `n` lines we get a lower bound on the code duplication factor of:

$$2n + 3/n + 1$$

For single line methods ( $n = 0$ ) the factor is 3 and for large methods ( $n \rightarrow \infty$ ) it converges to 2.

**Table 3.** Comparison of ScalaCT with LMS and hand optimized code.

Benchmark	ScalaCT	LMS	Optimized by Hand
pow			
min			
dot			
fft			

## 5.2 Performance of Generated Code

## 6. Limitations

- Interaction with type variables.
- Type variables.
- Type annotations and overloading and implicit search.
- Can not inherit from a compile time view.

## 7. Related Work

MetaOCaml [5, 34] is a staging extension of OCaml. It uses quotation to promote

MacroML [13]

Programming language Idris [3] try to solve this problem by allowing the `static` annotation on function parameters. Annotation `static` denotes that the term is statically known and that all operations on that term should be executed at compile-time. However, since `static` is placed on terms rather than types, it can mark only *whole terms* as static. This restricts the number of programs that can be expressed, e.g., we could not express that vectors in the signature of `dot` are static only in size. Finally, information about `static` terms can not be propagated through return types of functions so `static` in Idris is a partial evaluation construct, i.e., it hints that partial evaluation should be applied if function arguments are static.

Programming language D [1] also introduces the `static` keyword for function parameters.

[19]

Hybrid Partial Evaluation (HPE) [28] is a technique for partial evaluation that does not perform binding time analysis (similarly to online partial evaluators) but relies on the user provided annotation CT<sup>1</sup>. HPE implementations exist for both Java and Scala [29]. Although, CT is used for partial evaluation, it does not affect typing of user programs. Furthermore, behavior of CT in context of generics is not described. ScalaCT can be seen as statically typed version of hybrid partial evaluation with support for parametric polymorphism. Due to the support for parametric polymorphism ScalaCT can express compile-time data structures with dynamic data.

Jovanovic et al. [18] in the Yin-Yang framework solve the problem of code duplication by generating reification and

<sup>1</sup> Name `ct` in ScalaCT is inspired by hybrid partial evaluation.

code generation logic based on original Scala code. With their approach there is no code duplication for the supported language features. However, not all of Scala is supported and all generated terms are generated for the next stage, thus, staging is not possible there is no notion of staging.

Forge by Sujeeth et al. [31] uses a DSL to declare a specification of the libraries. Forge then generates both unannotated and annotated code based on the specification. Their language also supports generating staged code (comprised of terms different from multiple stages). Forge specification and code generation supports only a subset of Scala guided towards the Delite [4, 30] framework.

## 8. Conclusion

## References

- [1] The d programming language. <http://dlang.org/>.
- [2] Kenichi Asai. Binding-time analysis for both static and dynamic expressions. *New Generation Computing*, 20(1):27–51, 2002.
- [3] Edwin C Brady and Kevin Hammond. Scrapping your inefficient engine: using partial evaluation to improve domain-specific language implementation. In *International Conference on Functional Programming (ICFP)*, 2010.
- [4] K. J Brown, A. K Sujeeth, H. J Lee, T. Rompf, H. Chafi, M. Odersky, and K. Olukotun. A heterogeneous parallel framework for domain-specific languages. In *Parallel Architectures and Compilation Techniques (PACT)*, 2011.
- [5] Cristiano Calcagno, Walid Taha, Liwen Huang, and Xavier Leroy. Implementing multi-stage languages using ASTs, gensym, and reflection. In *Generative Programming and Component Engineering*, 2003.
- [6] Jacques Carette and Oleg Kiselyov. Multi-stage programming with functors and monads: Eliminating abstraction overhead from generic code. In *Generative Programming and Component Engineering (GPCE)*, 2005.
- [7] Krzysztof Czarnecki, John O'Donnell, Jörg Striegnitz, and Walid Taha. DSL implementation in MetaOCaml, Template Haskell, and C++. *Domain-Specific Program Generation*, 2004.
- [8] Olivier Danvy. *Type-directed partial evaluation*. Springer, 1999.
- [9] Rowan Davies. A temporal-logic approach to binding-time analysis. In *Symposium on Logic in Computer Science (LICS)*, 1996.
- [10] Rowan Davies and Frank Pfenning. A modal analysis of staged computation. In *Symposium on Principles of Programming Languages (POPL)*, 1996.
- [11] Zachary DeVito, James Hegarty, Alex Aiken, Pat Hanrahan, and Jan Vitek. Terra: a multi-stage language for high-performance computing. In *International Conference on Programming Language Design and Implementation (PLDI)*, 2013.
- [12] Yoshihiko Futamura. Partial evaluation of computation process—an approach to a compiler-compiler. *Higher-Order*

and *Symbolic Computation*, 12(4):381–391, 1999.

- [13] Steven E Ganz, Amr Sabry, and Walid Taha. Macros as multi-stage computations: type-safe, generative, binding macros in MacroML. In *International Conference on Functional Programming (ICFP)*, 2001.
- [14] Carsten K Gomard and Neil D Jones. A partial evaluator for the untyped lambda-calculus. *Journal of Functional Programming*, 1(01):21–69, 1991.
- [15] Cordelia V. Hall, Kevin Hammond, Simon L. Peyton Jones, and Philip L. Wadler. Type classes in Haskell. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 18(2):109–138, 1996.
- [16] Neil D. Jones, Carsten K. Gomard, and Peter Sestoft. *Partial Evaluation and Automatic Program Generation*. Prentice Hall, 1993.
- [17] Manohar Jonnalagedda, Thierry Coppey, Sandro Stucki, Tiark Rompf, and Martin Odersky. Staged parser combinators for efficient data processing. In *International Conference on Object Oriented Programming Systems Languages and Applications (OOPSLA)*, 2014.
- [18] V. Jovanovic, A. Shaikhha, S. Stucki, V. Nikolaev, Koch C., and M Odersky. Yin-Yang: Concealing the deep embedding of DSLs. In *International Conference on Generative Programming and Component Engineering (GPCE)*, 2014.
- [19] Anne-Françoise Le Meur, Julia L Lawall, and Charles Consel. Specialization scenarios: A pragmatic approach to declaring program specialization. *Higher-Order and Symbolic Computation*, 17(1-2):47–92, 2004.
- [20] Stefan Monnier and Zhong Shao. Inlining as staged computation. *Journal of Functional Programming*, 13(03):647–676, 2003.
- [21] F. Nielson and R. H. Nielson. Automatic binding time analysis for a typed  $\lambda$ -calculus. In *Symposium on Principles of Programming Languages (POPL)*, 1988.
- [22] Flemming Nielson and Hanne Riis Nielson. *Two-level functional languages*, volume 34. Cambridge University Press, 2005.
- [23] Martin Odersky and Konstantin Läufer. Putting type annotations to work. In *Symposium on Principles of Programming Languages (POPL)*, 1996.
- [24] Bruno C.d.S. Oliveira, Adriaan Moors, and Martin Odersky. Type classes as objects and implicits. In *Object-Oriented Programming Systems, Languages, and Applications (OOPSLA)*, 2010.
- [25] Tiark Rompf and Martin Odersky. Lightweight modular staging: a pragmatic approach to runtime code generation and compiled DSLs. *Communications of the ACM*, 55(6):121–130, 2012.
- [26] Tiark Rompf, Arvind K. Sujeeth, Nada Amin, Kevin J. Brown, Vojin Jovanović, HyoukJoong Lee, Manohar Jonnalagedda, Kunle Olukotun, and Martin Odersky. Optimizing data structures in high-level programs: New directions for extensible compilers based on staging. In *Symposium on Principles of Programming Languages (POPL)*, 2013.
- [27] Tiark Rompf, Arvind K. Sujeeth, Kevin J. Brown, HyoukJoong Lee, Hassan Chafi, Kunle Olukotun, and Martin Odersky. Project Lancet: Surgical precision JIT compilers. In *International Conference on Programming Language Design and Implementation (PLDI)*, 2013.
- [28] Amin Shali and William R. Cook. Hybrid partial evaluation. In *International Conference on Object Oriented Programming Systems Languages and Applications (OOPSLA)*, 2011.
- [29] Amanj Sherwany, Nosheen Zaza, and Nathaniel Nystrom. A refactoring library for scala compiler extensions. In *Compiler Construction (CC)*, pages 31–48, 2015.
- [30] Arvind Sujeeth, Tiark Rompf, Kevin Brown, HyoukJoong Lee, Hassan Chafi, Victoria Popic, Michael Wu, Aleksander Prokopec, Vojin Jovanović, Martin Odersky, and Kunle Olukotun. Composition and reuse with compiled domain-specific languages. In *European Conference on Object-Oriented Programming (ECOOP)*, 2013.
- [31] Arvind K Sujeeth, Austin Gibbons, Kevin J Brown, HyoukJoong Lee, Tiark Rompf, Martin Odersky, and Kunle Olukotun. Forge: generating a high performance DSL implementation from a declarative specification. In *International Conference on Generative Programming and Component Engineering (GPCE)*, 2013.
- [32] Eijiro Sumii and Naoki Kobayashi. A hybrid approach to on-line and offline partial evaluation. *Higher-Order and Symbolic Computation*, 14(2-3):101–142, 2001.
- [33] Walid Taha. A gentle introduction to multi-stage programming. In *Domain-Specific Program Generation (DSPG)*. 2004.
- [34] Walid Taha and Tim Sheard. Multi-stage programming with explicit annotations. In *Workshop on Partial Evaluation and Program Manipulation (PEPM)*, 1997.
- [35] Philip Wadler and Stephen Blott. How to make ad-hoc polymorphism less ad hoc. In *Symposium on Principles of Programming Languages (POPL)*. ACM, 1989.
- [36] Edwin Westbrook, Mathias Ricken, Jun Inoue, Yilong Yao, Tamer Abdelatif, and Walid Taha. Mint: Java multi-stage programming using weak separability. In *International Conference on Programming Language Design and Implementation (PLDI)*, 2010.