

Compile-Time Views: Predictable Type-Directed Partial Evaluation Without Code Duplication

No Author Given

No Institute Given

Abstract.

Keywords: Partial Evaluation

1 Introduction

Partial evaluation [5] is an optimization technique that identifies *statically known* program parts and pre-computes them at compile time. The compile-time computation yields a *residual program* that does not contain the, previously identified, statically known parts of the program. Partial evaluation has been intensively studied and successfully applied for: removing abstraction overheads in high-level programs [2, 9], domain-specific languages [1, 6], and converting language interpreters into compilers [4, 10, 13]. Applying partial evaluation in these domains often improves program performance by several orders of magnitude [11, 1].

Unlike other compiler optimizations partial evaluation is not *safe*: it might lead to *code explosion* and might not *terminate*. Due to compile-time execution, computing **folds** and loops over data structures of static size can produce arbitrarily large residual programs. Furthermore, in a Turing-complete language assuring termination of partial-evaluation is undecidable. **[TODO: cite]**

Automatically assuring safety of partial evaluators necessarily leads to lack of *predictability*. To illustrate, let us define a function **dot** for computing a dot-product of two vectors that contain numeric values¹:

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

When **dot** is called with vectors of static size (*e.g.* `dot(Vector(2, 4), Vector(1, 10))`) the abstraction overhead of **zip** and **foldLeft** can be completely removed. However, the partial evaluator must apply extensive analysis to conclude that vectors are of static size and that this information can be later used to unroll the recursion inside **foldLeft**. Even if the analysis is successful the evaluator must be conservative about unrolling the **foldLeft**. The vector sizes, and thus the produced code, can unacceptably large in a general case. What if we know

that vector sizes are relatively small and we would like to predictably unroll `dot` into a flat sum of products?

Lack of predictability and danger of code explosion are the reason that successful partial evaluators [1, 12, 9, 13, 7] are programmer controlled. We categorize the existing solutions in three categories (for further discussion *c.f.* §7):

- Programming languages Idris and D provide allow placing the `static` annotation on function arguments. Since `static` is placed on terms, it denotes that the *whole term* is 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 partially static.
- Type-directed partial evaluation [3] and Lightweight Modular Staging (LMS) [9] use types to communicate the programmer’s intent about partial evaluation. By changing the types of parameters to be (*e.g.* `Vector[Rep[T]]`) these approaches can express that parameter vectors are statically known. However, they still require existence of two data structures (*e.g.* `Rep[Vector]` and `Vector`). This fosters costly and hardly maintainable code duplication.
- MetaOCaml [12] places terms in, possibly nested, quotes. Depth of the term in the quotes denotes the stage of the computation where it will be executed. In MetaOCaml we can express the `dot` function, but we have to modify the code of the `dot` function which might not be desirable.

Ideally, a programmer would with a minimal number of annotations be able to: *i)* require that input vectors are of statically known size but polymorphic in their elements, *ii)* without modifying the terms require that all operations on vector arguments are further partially evaluated, *iii)* allow elements of vectors to be generic, and *iv)* reuse the existing implementation of the `Vector` data structure.

The main idea of this paper is to provide a statically typed *compile-time view* of existing data types. The compile-time view makes all operations and non-generic fields partially evaluated on a type. The compile-time view allows programmers to define a single definition of a type. Then the existing types can be promoted to their compile-time duals with the `@ct` annotation at the type level, and with the `ct` function on the term level. Consequently, due to the integration with the type system, the control over partial evaluation is fine-grained and polymorphic and term level promotions obviate code duplication for static data structures.

With our partial evaluator, to require that vectors `v1` and `v2` are static and to partially evaluate the function, a programmer would need to make a simple modification of the `dot` signature:

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

This, in effect, requires that only vector arguments (not their elements) are statically known and that all operations on vector arguments will be executed at compile time (partially evaluated). Since, values are polymorphic the result of the function will either be a dynamic value or a compile-time value. Residual programs of `dot` application for different arguments:

```
// [el1, el2, el3, el4] are dynamic
dot(ct(Vector)(el1, el2), ct(Vector)(el3, el4))
  ⇨ el1 * el3 + el2 * el4

dot(ct(Vector)(2, 4), ct(Vector)(1, 10))
  ⇨ 2 * 1 + 4 * 10

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

In this paper we make the following contributions to the state-of-the-art:

- By introducing the $F_{i<}$ calculus (§4) that in a fine-grained way captures the user’s intent about partial evaluation. The calculus is based on $F_{<}$ with lazy records which makes it suitable for representing modern multi-paradigm languages with object oriented features. Finally, we formally define a partial evaluator for $F_{i<}$.
- By providing a *translation scheme* from data types in object oriented languages (polymorphic classes and methods) into their dual compile-time views in the $F_{i<}$ calculus (§5).
- 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 [8].

We have implemented a partial evaluator according to the translation scheme (§5) from object oriented features of Scala to the $F_{i<}$ calculus. The partial is implemented for Scala and open-sourced (<https://github.com/scala-inline/>). It has a minimal Scala interface (§2) based on type annotations. We have evaluated the performance gains and the validity of the partial evaluator on all case studies (§3) and compared them to LMS. In all benchmarks our evaluator gives significant performance gains compared to original programs and performs equivalently to LMS.

2 The Partial Evaluator for Scala

We have implemented a prototype partial evaluator, formally defined in §??, and according to the $F_{i<}$ calculus (formally define in §4). The partial evaluator is a compiler plugin that executes in a phase after the Scala type checker. The plugin starts with pre-typed Scala programs and uses a type annotations **[TODO: cite]** to track and verify information about the binding-time of terms.

To the user, the partial evaluator exposes a minimal interface (Figure 2) with annotations `inline` and `ct` and the `ct` function.

Annotation `ct` is used at the type level and denotes that one expects a compile-time view of a type. The annotation is integrated in the Scala’s type system and, therefore, can be arbitrarily nested in different variants of types. Table 2 shows how the `@ct` annotation can be placed on types and how it,

```

package object scalainline {

  final class ct extends StaticAnnotation
  final class inline extends StaticAnnotation

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

}

```

Fig. 1. Interface of the Scala partial evaluator.

due to the translation to the compile-time views (Figure ??), changes method signatures.

Table 1. Types and corresponding method signatures after the translation to the compile-time view.

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 => U)@ct): Vector[U]@ct</code> <code>length: Int@ct</code>
<code>Vector[Int@ct]@ct</code>	<code>map[U](f: (Int@ct => U)@ct): Vector[U]@ct</code>
<code>Map[Int@ct, Int]@ct</code>	<code>get(key: Int@ct): Option[Int]@ct</code>

In Table 2, `Int@ct` is a non-polymorphic type and therefore according to the translation to the compile-time view (13) all arguments of all methods will be executed at compile-time. On the other hand, `Vector[Int]@ct` will have all arguments 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` is has all parts executed at compile-time. However, the return type of the function `map` can still be a compile-time view - due to the type parameter `U`.

Functions `ct` and `inline` is used at the term level for promoting Scala objects and functions into their compile-time views. Without `ct` we would not be able to instantiate compile-time views of the types. Table 2 shows how different types of terms are promoted to their compile-time views.

[TODO: footnote about Scala objects] [TODO: static promotion of lambdas] Function `ct` can be applied to objects (*e.g.* `Vector`) to provide a compile-time view over their methods. When those objects have generic parameters, `ct` be used to promote the arguments, and thus, the result types of these functions. When applied, on functions `ct` promotes the compile-time view as well as its arguments and the return type. **[TODO: inline]**

Annotation `inline` can be used only on methods and functions. This function uses partial evaluation to achieve inlining**[TODO: cite]**. This is not the

Table 2. Types and corresponding method signatures after the translation to the compile-time view.

Promoted Term	Term's Promoted 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>new (Cons@ct)(1, Nil)</code>	<code>: Cons[Int]@ct</code>
<code>new (Cons@ct)(ct(1), ct(Nil))</code>	<code>: Cons[Int@ct]@ct</code>
<code>ct((x: Int) => x)</code>	<code>: (Int@ct => Int@ct)@ct</code>
<code>inline((x: Int) => x)</code>	<code>: (Int => Int)@ct</code>

first time that inlining is achieved through partial evaluation[**TODO: cite**], however, partial evaluation is trivially added to the system. It directly corresponds to adding `inline` from $F_{i<}$ in front of the function or method definition.

2.1 Interaction with the Scala Language

3 Case Studies

In this section we present selected use-cases for compile-time views that demonstrate the core functionality. We start with a canonical example of the 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-classes can be removed (§3.5).

3.1 Inlining Expressed Through Partial-Evaluation

3.2 Recursion

The canonical example in partial evaluation is the computation of the integer power function:

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

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

With compile-time views making `pow` partially evaluated requires adding two annotations:

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

`@inline` denotes that the `pow` function itself must be inlined at application and `@ct` requires that the `exp` argument is a compile-time view of `Int`. The application of the function `pow` with a constant exponent will produce:

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

Here, in the function application, constant 4 is promoted to `ct` by the automatic conversions. **[TODO: ref]**

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

```
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 }
```

Fig. 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`.

3.4 Removing Abstraction Overhead of Type-Classes

[TODO: not-sure how to achieve this!] [TODO: cite] Type-classes are omnipresent in everyday programming as they provide allow abstraction over generic parameters (*e.g.* Numeric abstracts over numeric values). Unfortunately, type-classes are a source of abstraction overheads during execution**[TODO: cite]**. Type-classes are in most of the cases statically known. Ideally, we would be able to deterministically remove abstraction overheads of type classes.

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

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

class DoubleNumeric[T <: Double] extends Numeric[Double] {
  @inline def plus(x: T, y: T): T = x + y
  @inline def times(x: T, y: T): T = x * y
  @inline def zero: T = 0.0
}
```

Fig. 3. Function for computing the non-negative power of a real number.

3.5 Dot Product

- Explain the removal of type classes together with inline. Explain how type classes are @i? and how they will completely evaluate if they are passed a static value.
- Comparison to other approaches.

4 The $F_{i<}$ Calculus

$t ::=$	Terms:	$S, T, U ::=$	Types:
x, y	identifier	$iS \Rightarrow jT$	function type
$(x : iT) \Rightarrow t$	function	$\{x : iS\}$	record type
$t(t)$	application	$[X <: iS] \Rightarrow jT$	universal type
$\{x = t\}$	record	Any	top type
$t.x$	selection	$iT, jT, kT, lT ::=$	Binding-Time Types:
$[X <: iT] \Rightarrow t$	type abstraction	X	type identifier
$t[iT]$	type application	$T, dynamic\ T$	dynamic type
$inline\ t$	inline view	$static\ T$	static type
$v ::=$	Values:	$inline\ T$	inline type
$x \Rightarrow t$	function value	$\Gamma ::=$	Contexts:
$\{x = t\}$	record value	\emptyset	empty context
		$\Gamma, x : iT$	term binding
		$\Gamma, X <: iT$	type binding

Fig. 4. Syntax of $F_{i<}$.

We formalize the essence of our inlining system in a minimalistic calculus based on $F_{<}$ with lazy records. To accommodate predictable partial evaluation we introduce binding-time annotations into the type system as first-class types that represent three kinds of bindings:

1. **Dynamic binding.** These are the types which express computation at run-time. All types written in the end user code are considered to be dynamic by default if no other binding-time annotation is given.
2. **Static binding.** Values of static terms can be computed at compile-time (*e.g.* constant expressions) but are still evaluated at runtime by default. All language literals are static by default.
3. **Inline binding.** And finally the types that correspond to terms that are hinted to be computed at compile-time whenever possible.

4.1 Composition

An interesting consequence of encoding of binding times as first-class types is ability to represent values which are partially static and partially dynamic.

For example lets have a look at simple record that describes a complex number with two possible representations encoded through *isPolar* flag:

$$complex : static \{isPolar : static\ Boolean, a : Double, b : Double\} \in \Gamma$$

This type is constructed out of a number of components with varying binding times. Representation encoding is known in advance and is static according to the signature. Coordinates a and b do not have any binding-time annotation meaning that they are dynamic.

Given this binding to *complex* in our environment Γ we can use *inline* to obtain a compile-time view to evaluate access to *isPolar* field at compile-time:

inline complex.isPolar : inline Boolean

Any statically known expression can be promoted via *inline*. Selection of dynamic fields on the other hand will return dynamic values despite the fact that record is statically known. In practice this can be used to specialize a particular execution path in the application to a particular representation by selectively inlining statically known parts.

Once you have inline view of the term it's also possible to demote it back to runtime evaluation through *dynamic* view.

Not all type and binding time combinations are correct though. We restrict types to disallow nesting of more specific binding times into less specific ones.

$$\begin{array}{ll}
 \text{wff } iAny & \text{(W-ANY)} \\
 \frac{i <: j \quad i <: k \quad \text{wff } jT_1 \quad \text{wff } kT_2}{\text{wff } i(jT_1 \Rightarrow kT_2)} & \text{(W-ABS)} \\
 \frac{i <: j \quad i <: k \quad \text{wff } jS \quad \text{wff } kT}{\text{wff } i([X <: jS] \Rightarrow kT)} & \text{(W-TABS)} \\
 \frac{\forall j. \quad i <: j \quad \text{wff } jT}{\text{wff } i\{x : jT\}} & \text{(W-REC)}
 \end{array}$$

Fig. 5. Well formed types wff iT

This restriction allows us to reject programs that have inconsistent annotations. For example the following function has incorrectly annotated parameter binding time:

$$(x : \text{inline Int}) \Rightarrow x + 1$$

This is inconsistent because the body of the function might not be evaluated at compile-time (as the function is not inline.) As described in (W-ABS) functions may only have parameters that are at most as specific as the function binding-time. In our example this doesn't hold as *inline* is more specific than implicit *static* annotation on function literal.

4.2 Subtyping

Another notable feature of our binding-time analysis system is deep integration with subtyping. We believe that such integration is crucial for an object-oriented language that wants to incorporate partial evaluation.

At core of the subtyping relation we have a subtyping relation on binding-time information with *dynamic* as top binding-time.

$$\begin{array}{ll} i <: \textit{dynamic} & (\text{I-DYNAMIC}) \\ \textit{static} <: \textit{static} & (\text{I-STATIC1}) \end{array} \qquad \begin{array}{ll} \textit{inline} <: \textit{static} & (\text{I-STATIC2}) \\ \textit{inline} <: \textit{inline} & (\text{I-INLINE}) \end{array}$$

Fig. 6. Binding-time subtyping.

We proceed by threading binding time information throughout regular $F_{<}$ subtyping rules augmented with standard record types.

$$\begin{array}{ll} \Gamma \vdash iS <: \textit{Any} & (\text{S-TOP}) \\ \Gamma \vdash iS <: iS & (\text{S-REFL}) \\ \frac{\Gamma \vdash iS <: jU \quad \Gamma \vdash jU <: kT}{\Gamma \vdash iS <: kT} & (\text{S-TRANS}) \\ \frac{i <: j \quad \Gamma \vdash S <: T}{\Gamma \vdash iS <: jT} & (\text{S-INLINE}) \\ \frac{X <: iT \in \Gamma}{\Gamma \vdash X <: iT} & (\text{S-TVAR}) \\ \frac{\{x_p : i_p T_p^{p \in 1..n+m}\} <: \{x_p : i_p T_p^{p \in 1..n}\}}{\Gamma \vdash kT_1 <: iS_1 \quad \Gamma \vdash jS_2 <: lT_2} & (\text{S-WIDTH}) \\ \frac{\Gamma \vdash kT_1 <: iS_1 \quad \Gamma \vdash jS_2 <: lT_2}{\Gamma \vdash iS_1 \Rightarrow jS_2 <: kT_1 \Rightarrow lT_2} & (\text{S-ARROW}) \\ \frac{\forall p \in 1..n. i_p S_p <: j_p T_p}{\{x_p : i_p S_p^{p \in 1..n}\} <: \{x_p : j_p T_p^{p \in 1..n}\}} & (\text{S-DEPTH}) \\ \frac{\Gamma, X <: iU_1 \vdash jS_2 <: kT_2}{\Gamma \vdash [X <: iU_1] \Rightarrow jS_2 <: [X <: iU_1] \Rightarrow kT_2} & (\text{S-ALL}) \\ \frac{\{x_p : i_p S_p^{p \in 1..n}\} \text{ is permutation of } \{y_p : j_p T_p^{p \in 1..n}\}}{\{x_p : i_p S_p^{p \in 1..n}\} <: \{y_p : j_p T_p^{p \in 1..n}\}} & (\text{S-PERM}) \end{array}$$

Fig. 7. Subtyping.

Integration between binding-time subtyping and subtyping on regular types is expressed through (S-INLINE) rule that merges the two into one coherent relation on binding-time types.

4.3 Generics

Crucial consequence of our design choices made in the system manifests in ability to use regular generics as means to abstract over binding-time without any additional language constructs.

For example given a generic identity function:

$$\textit{identity} : \textit{static} ([X <: \textit{Any}] \Rightarrow \textit{static} (X \Rightarrow X)) \in \Gamma$$

We can instantiate it to both in static and dynamic contexts through corresponding type application:

$$\begin{aligned} \text{identity}[\text{static } Int] &: \text{static } (\text{static } Int \Rightarrow \text{static } Int) \\ \text{identity}[Int] &: \text{static } (Int \Rightarrow Int) \end{aligned} \quad (1)$$

In practice this allows us to write code that is polymorphic in the binding time without any code duplication which is quite common in other partial evaluation systems.

This is possible due to the fact that we've integrated binding time information into types and augmented subtyping relation with subtyping

4.4 Typing

To enforce well-formedness of types in a context of partial evaluation we customize standard typing rules with additional constraints with respect to binding time.

$$\begin{aligned} & \frac{x : iT \in \Gamma}{\Gamma \vdash x : iT} & (\text{T-IDENT}) \\ & \frac{\forall t. \quad \Gamma \vdash t : jT \quad \text{wff } i\{x : jT\}}{\Gamma \vdash i\{x = t\} : i\{x : jT\}} & (\text{T-REC}) \\ & \frac{\Gamma \vdash t_1 : i(jT_1 \Rightarrow kT_2) \quad \Gamma \vdash t_2 : jT_1}{\Gamma \vdash t_1(t_2) : kT_2} & (\text{T-APP}) \\ & \frac{\Gamma \vdash t : i\{x = jT_1, y = kT_2\}}{\Gamma \vdash t.x : jT_1} & (\text{T-SEL}) \\ & \frac{t \text{ is not literal} \quad \Gamma \vdash t : \text{static } T}{\Gamma \vdash \text{inline } t : \text{inline } T} & (\text{T-INLINE}) \\ & \frac{t \text{ is not literal} \quad \Gamma \vdash t : iT}{\Gamma \vdash \text{dynamic } t : \text{dynamic } T} & (\text{T-DYNAMIC}) \\ & \frac{\Gamma \vdash t : iS \quad \Gamma \vdash iS <: jT}{\Gamma \vdash t : jT} & (\text{T-SUB}) \\ & \frac{\Gamma, x : jT_1 \vdash t : kT_2 \quad \text{wff } i(jT_1 \Rightarrow kT_2)}{\Gamma \vdash i((x : jT_1) \Rightarrow t) : i(jT_1 \Rightarrow kT_2)} & (\text{T-FUNC}) \\ & \frac{\Gamma, X <: jT_1 \vdash t_2 : kT_2 \quad \text{wff } i([X <: jT_1] \Rightarrow kT_2)}{\Gamma \vdash i([X <: jT_1] \Rightarrow t_2) : i([X <: jT_1] \Rightarrow kT_2)} & (\text{T-TABS}) \\ & \frac{\Gamma \vdash t_1 : i([X <: jT_{11}] \Rightarrow kT_{12}) \quad \Gamma \vdash lT_2 <: jT_{11} \quad \Gamma \vdash i <: l}{\Gamma \vdash t_1[lT_2] : [X \mapsto lT_2]kT_{12}} & (\text{T-TAPP}) \end{aligned}$$

Fig. 8. Typing.

The most significant changes lie in:

- Additional checks in literal typing that ensure that constructed values correspond to well-formed types (T-FUNC, T-REC, T-TABS). To do this we typecheck literals together with possible binding-time term that might enclose it.
- New typing rules for binding-time views (T-INLINE, T-DYNAMIC). These rules only cover non-literal terms as composition of binding-time view and literal itself is handled in corresponding typing rule for given literal.

4.5 Partial Evaluation

In order to simplify partial evaluation rules we erase all of the type information before partial evaluation. This means that all functions become function values, type abstraction and application are complete eliminated.

$$\begin{array}{c}
\frac{t \rightsquigarrow t'}{x \Rightarrow t \rightsquigarrow x \Rightarrow t'} \quad (\text{PE-FUNC}) \\
\frac{\bar{t} \rightsquigarrow \bar{t}'}{\{\bar{x} = \bar{t}\} \rightsquigarrow \{\bar{x} = \bar{t}'\}} \quad (\text{PE-REC}) \\
\frac{t_1 \rightsquigarrow t'_1 \quad t'_1 \neq \text{inline } x \Rightarrow t \quad t_2 \rightsquigarrow t'_2}{t_1(t_2) \rightsquigarrow t'_1(t'_2)} \quad (\text{PE-APP}) \\
\frac{t_1 \rightsquigarrow \text{inline } x \Rightarrow t \quad t_2 \rightsquigarrow t'_2 \quad [x \mapsto t'_2]t \rightsquigarrow t'}{t_1(t_2) \rightsquigarrow t'} \quad (\text{PE-IAPP}) \\
\frac{t \rightsquigarrow t' \quad t' \neq \text{inline } x \Rightarrow t}{t.x \rightsquigarrow t'.x} \quad (\text{PE-SEL}) \\
\frac{t \rightsquigarrow \text{inline } \{x = t_x, \bar{y} = \bar{t}_y\} \quad t_x \rightsquigarrow t'_x}{t.x \rightsquigarrow t'_x} \quad (\text{PE-ISEL}) \\
\frac{t \text{ is not literal} \quad t \rightsquigarrow t' \quad t' \Downarrow v}{\text{inline } t \rightsquigarrow \text{inline } v} \quad (\text{PE-INLINE})
\end{array}$$

Fig. 9. Partial evaluation $t \rightsquigarrow t'$

4.6 Evaluation

Once partial evaluation is complete we strip all binding-time terms and use regular untyped lambda calculus evaluation rules extended with lazy records.

$$\begin{array}{c}
\frac{v \Downarrow v}{t_1 \Downarrow x \Rightarrow t \quad t_2 \Downarrow v \quad [x \mapsto v]t \Downarrow v'} \quad (\text{E-APP}) \\
\frac{t \Downarrow \{x = t_x, \bar{y} = \bar{t}_y\} \quad t_x \Downarrow v}{t.x \Downarrow v} \quad (\text{E-SEL})
\end{array}$$

Fig. 10. Evaluation $t \Downarrow v$

4.7 Conjectures

1. Progress.
2. Preservation.
3. Static terms are closed over statically bound variables.
4. Inline terms will be replaced with canonical value of corresponding type after partial evaluation.

5 Integrating $F_{i<}$ with Object Oriented Languages

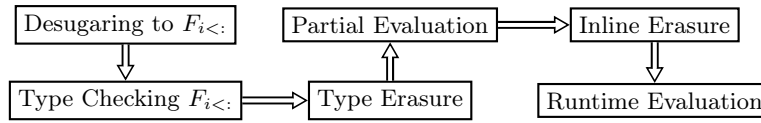


Fig. 11. Compilation pipeline.

The $F_{i<}$ calculus §4 captures the essence of user-controlled predictable partial-evaluation. In practice, though, it is fairly low level and it is not obvious how to define *classes* and methods from in modern multi-paradigm programming languages. Furthermore, $F_{i<}$ requires an inconveniently large number of `inline` calls in method invocations. In this section we a scheme for translating classes into $F_{i<}$ (§5.1), show how to provide compile time views of classes and *methods*§??, and formalize convenient implicit conversions for the calculus §5.3.

Furthermore, rules of $F_{i<}$ do not support effect-full computations and each `inline` term is trivially converted to a dynamic term after erasure. In case of languages that do support mutable state and side-effects this needs to be treated specially. For simplicity, we omit side-effects from our discussion and assume that all partially evaluated code is side-effect free and that each `inline` term can be converted to dynamic code.

5.1 Desugaring Object Oriented Constructs to $F_{i<}$

5.2 Compile-Time View of the Terms

5.3 Implicit Conversions

According to $F_{i<}$ rules if method signatures contain compile-time views of a type the corresponding arguments in method application would always have to be promoted to `inline`. In practice this is not convenient as it requires an inconveniently large number of annotations. Partial evaluation is an optimization, and as such, it should not affect user code - users should not be aware of the internal operation of the library.

$\llbracket \text{let } x : T_x = t_x \text{ in } t \rrbracket = ((x : T_x) \Rightarrow t)(t_x)$
 $\llbracket \text{let type } T_1 = T_2 \text{ in } t \rrbracket = ([T_1 <: T_2] \Rightarrow t)[T_2]$
 $\llbracket \text{let class } C[A](x : T_x) \{ \text{def } f[B](y : T_y) = t_f \} \text{ in } t \rrbracket =$
 $\text{let type } C = [A] \Rightarrow \text{inline } \{ \text{fields} : \{ x : T_x \}, \text{methods} : \text{inline } \{ f : [B] \Rightarrow T_y \Rightarrow T_f \} \} \text{ in}$
 $\text{let } C : [A] \Rightarrow \text{inline } ((t_x : T_x) \Rightarrow C[A]) = [A] \Rightarrow \text{inline } ((t_x : T_x) \Rightarrow$
 $\text{inline } \{ \text{fields} = \{ x = t_x \}, \text{methods} = \text{inline } \{ f = [B] \Rightarrow (y : T_y) \Rightarrow t_f \} \}) \text{ in } t$

Fig. 12. Desugaring of classes into $F_{i<}$.

$$\begin{array}{c}
\frac{\Pi \vdash T \in \Pi}{\Pi \vdash iT \rightsquigarrow iT} \text{ (CT-TVAR)} \qquad \frac{\Pi \vdash T \notin \Pi}{\Pi \vdash iT \rightsquigarrow \text{inline } T} \text{ (CT-T-VAR)} \\
\\
\frac{\Pi \vdash t \rightsquigarrow t'}{\Pi \vdash i\{x = t\} \rightsquigarrow \text{inline } \{x = t'\}} \text{ (CT-REC)} \\
\frac{\Pi \vdash iT \rightsquigarrow jT}{\Pi \vdash \{x : iT\} \rightsquigarrow \text{inline } \{x : jT\}} \text{ (CT-T-REC)} \\
\frac{\Pi \vdash iT \rightsquigarrow jT \quad \Pi \vdash kS \rightsquigarrow lS}{\Pi \vdash iT \Rightarrow kS \rightsquigarrow jT \Rightarrow lS} \text{ (CT-T-ARROW)} \\
\frac{\Pi \vdash jT \rightsquigarrow kT}{\Pi \vdash [X <: iS] \Rightarrow jT \rightsquigarrow [X <: iS] \Rightarrow kT} \text{ (CT-T-UNIV)} \\
\frac{\Pi \vdash t \rightsquigarrow t' \quad \Pi \vdash iT \rightsquigarrow jT}{\Pi \vdash i(x : iT) \Rightarrow t \rightsquigarrow \text{inline } (x : jT) \Rightarrow t'} \text{ (CT-FUNC)} \\
\frac{\Pi, X \vdash t \rightsquigarrow t'}{\Pi \vdash i([X <: jT_1] \Rightarrow t) \rightsquigarrow \text{inline } ([X <: jT_1] \Rightarrow t')} \text{ (CT-TABS)} \\
\frac{\Pi \vdash t \rightsquigarrow t' \quad \Pi \vdash iT \rightsquigarrow jT}{\Pi \vdash t[iT] \rightsquigarrow t'[jT]} \text{ (CT-TAPP)}
\end{array}$$

Fig. 13. Translation of a type abstractions, function, and record values into a compile-time view. The translation is used for promoting types into their compile time versions.

To address this issue we introduce implicit conversions from all language literals, and direct class constructor calls of non-inline type into their compile-time views. For example, for a factorial function

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

we will not require annotations on literals 0, and 1. Furthermore, the function can be invoked without promoting the literal 5 into it's compile-time view:

```
fact(5)
  ↪ 120
```

6 Evaluation

6.1 Reduction of Code Duplication

6.2 Performance Comparison

Table 3. Performance comparison with LMS and hand optimized code.

Benchmark	Hand Optimized	LMS	Scala Inline
<code>pow</code>			
<code>min</code>			
<code>dot</code>			
<code>fft</code>			

7 Related Work

8 Conclusion

References

1. 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.
2. 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.
3. Olivier Danvy. *Type-directed partial evaluation*. Springer, 1999.
4. Yoshihiko Futamura. Partial evaluation of computation process—an approach to a compiler-compiler. *Higher-Order and Symbolic Computation*, 12(4):381–391, 1999.
5. Neil D. Jones, Carsten K. Gomard, and Peter Sestoft. *Partial Evaluation and Automatic Program Generation*. Prentice Hall, 1993.

6. 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.
7. 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.
8. Bruno CdS Oliveira, Adriaan Moors, and Martin Odersky. Type classes as objects and implicits. In *ACM Sigplan Notices*, volume 45, pages 341–360, 2010.
9. 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, June 2012.
10. 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.
11. Amin Shali and William R. Cook. Hybrid partial evaluation. In *International Conference on Object Oriented Programming Systems Languages and Applications (OOPSLA)*, 2011.
12. Walid Taha and Tim Sheard. Multi-stage programming with explicit annotations. In *Workshop on Partial Evaluation and Program Manipulation (PEPM)*, 1997.
13. Thomas Würthinger, Christian Wimmer, Andreas Wöß, Lukas Stadler, Gilles Duboscq, Christian Humer, Gregor Richards, Doug Simon, and Mario Wolczko. One VM to rule them all. In *Symposium on New Ideas, New Paradigms, and Reflections on Programming & Software (Onward!)*, 2013.