

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, (c1, cr)) =>  
    prod + c1 * 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 (§2) 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 (§3).
- By demonstrating the usefulness of compile-time views in four case studies (§4): 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 (§3) 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 (§5) based on type annotations. We have evaluated the performance gains and the validity of the partial evaluator on all case studies (§4) 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 $F_{i<}$ Calculus

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

Fig. 1. Syntax of $F_{i<}$.

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

1. **Dynamic binding.** Corresponds to terms that are expected to be evaluated at runtime.
2. **Inline binding.** Corresponds to terms that must be evaluated at compile-time.

To simplify judgments in our formalization we use concise iT syntax to abstract over binding-times of types. Here i signifies the bit of information that says if type is inline or not and T carries underlying type that is being annotated. So for example in *inline Any* we get $i = inline$ and $T = Any$.

Similarly we abstract over binding time of terms through it notation that has analogous to the one we use for types.

2.1 Well-formed types

Even though binding-time information is represented as types, not all of the possible combinations of types and binding-times is correct. We restrict types to disallow nesting of more specific binding times into less specific ones.

$$\begin{array}{ll}
& \text{wff } iAny & (\text{W-ANY}) \\
\frac{i \leq j \quad i \leq k \quad \text{wff } jT_1 \quad \text{wff } kT_2}{\text{wff } i(jT_1 \Rightarrow kT_2)} & (\text{W-ABS}) \\
\frac{i \leq k \quad \text{wff } [X \mapsto iS]kT}{\text{wff } i([X <: iS] \Rightarrow kT)} & (\text{W-TABS}) \\
\frac{i \leq \bar{j} \quad \text{wff } j\bar{T}}{\text{wff } i\{x : j\bar{T}\}} & (\text{W-REC})
\end{array}$$

Fig. 2. Well-formed types wff iT .

We represent notion of more specific binding-times through a simple partial order on binding time annotations.

$$\begin{array}{l}
dynamic \leq dynamic \\
inline \leq dynamic \\
inline \leq inline
\end{array}$$

Fig. 3. Partial order on binding-time $i \leq j$

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

$$(x : inline \ Int) \Rightarrow x + 1$$

This is inconsistent because a dynamic function may not have any non-dynamic parameters. As described in W-ABS functions may only have parameters that are at most as specific as function binding-time. In our example this doesn't hold as *inline* is more specific than *dynamic*.

2.2 Subtyping

$F_{i<}$: integrates binding-time annotation into subtyping relation on regular types by threading inlining information throughout all of the standard subtyping rules.

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

Fig. 4. Subtyping $\Gamma \vdash T_1 <: T_2$.

Apart from that we also introduce subtyping on binding-time types.

$$\begin{array}{c}
\frac{X <: iT \in \Gamma}{\Gamma \vdash X <: iT} \quad (S\text{-TVAR}) \\
\frac{i = j \quad \Gamma \vdash S <: T}{\Gamma \vdash iS <: jT} \quad (S\text{-INLINE})
\end{array}$$

Fig. 5. Subtyping of binding-time types $\Gamma \vdash iT_1 <: jT_2$.

Two binding-time types are subtypes if their underlying types are subtypes and if they have the same binding time.

2.3 Type polymorphism

Our system retains traditional type abstraction means inherited from $F_{<}$. We extend it to accomodate encoding of binding-times into types. This allows us to specify binding type of the abstracted generic type:

$$[T <: \textit{dynamic Any}] \Rightarrow (x : T) \Rightarrow x$$

For this particular identity function we need to restrict subset of all admissible types to only allow *dynamic* ones. Passing an *inline* type would not make sense as the resulting type would have not been well-formed.

2.4 Typing

$\frac{x : iT \in I}{\Gamma \vdash x : iT}$	(T-IDENT)
$\frac{\Gamma \vdash \bar{t} : \overline{jT} \quad \text{wff } i\{\overline{x : jT}\}}{\Gamma \vdash i\{\overline{x = t}\} : i\{\overline{x : jT}\}}$	(T-REC)
$\frac{\Gamma \vdash t_1 : i(jT_1 \Rightarrow kT_2) \quad \Gamma \vdash t_2 : jT_1}{\Gamma \vdash t_1(t_2) : kT_2}$	(T-APP)
$\frac{\Gamma \vdash t : i\{x = jT_1, \overline{y = kT_2}\}}{\Gamma \vdash t.x : jT_1}$	(T-SEL)
$\frac{\Gamma \vdash t : iS \quad \Gamma \vdash iS <: jT}{\Gamma \vdash t : jT}$	(T-SUB)
$\frac{\Gamma, x : jT_1 \vdash t : kT_2 \quad \text{wff } i(jT_1 \Rightarrow kT_2)}{\Gamma \vdash i((x : jT_1) \Rightarrow t) : i(jT_1 \Rightarrow kT_2)}$	(T-FUNC)
$\frac{\Gamma, X <: jT_1 \vdash t_2 : kT_2 \quad \text{wff } i([X <: jT_1] \Rightarrow kT_2)}{\Gamma \vdash i([X <: jT_1] \Rightarrow t_2) : i([X <: T_1] \Rightarrow kT_2)}$	(T-TABS)
$\frac{\Gamma \vdash t : i([X <: T_1] \Rightarrow kT_2) \quad \Gamma \vdash T <: T_1}{\Gamma \vdash t[T] : [X \mapsto T]kT_2}$	(T-TAPP)
$\frac{\Gamma \vdash t : \text{inline } T}{\Gamma \vdash \text{dynamic } t : \text{dynamic } T}$	(T-DYNAMIC)

Fig. 6. Typing $\Gamma \vdash t : iT$.

Similarly to the changes made to the subtyping relation we thread binding-time information throughout typing relation. Apart from that we also ensure that all literals produces by the user have well-formed types.

2.5 Partial Evaluation

$$\begin{array}{c}
\frac{t \rightsquigarrow t'}{(x : iT) \Rightarrow t \rightsquigarrow (x : T) \Rightarrow t'} \quad (\text{PE-FUNC}) \\
\frac{\bar{t} \rightsquigarrow \bar{t}'}{\{x = t\} \rightsquigarrow \{x = t'\}} \quad (\text{PE-REC}) \\
\frac{t \rightsquigarrow t'}{[X <: iT] \Rightarrow t \rightsquigarrow [X <: iT] \Rightarrow t'} \quad (\text{PE-TABS}) \\
\frac{t_1 \rightsquigarrow t'_1 \quad t_1 \neq \text{inline } t_3 \quad t_2 \rightsquigarrow t'_2}{t_1(t_2) \rightsquigarrow t'_1(t'_2)} \quad (\text{PE-APP}) \\
\frac{t \rightsquigarrow t' \quad t' \neq \text{inline } t_3}{t.x \rightsquigarrow t'.x} \quad (\text{PE-SEL}) \\
\frac{t_1 \rightsquigarrow t'_1 \quad t'_1 \neq \text{inline } t_3}{t_1[T_1] \rightsquigarrow t'_2} \quad (\text{PE-TAPP}) \\
\frac{t_1 \rightsquigarrow \text{inline } (x : iT) \Rightarrow t \quad t_2 \rightsquigarrow t'_2 \quad [x \mapsto t'_2]t \rightsquigarrow t'}{t_1(t_2) \rightsquigarrow t'} \quad (\text{PE-INLINEAPP}) \\
\frac{t \rightsquigarrow \text{inline } \{x = t_x, \overline{y = t_y}\} \quad t_x \rightsquigarrow t'_x}{t.x \rightsquigarrow t'_x} \quad (\text{PE-INLINESEL}) \\
\frac{t_1 \rightsquigarrow \text{inline } [X <: iT_2] \Rightarrow t_2 \quad [X \mapsto iT_1]t_2 \rightsquigarrow t'_2}{t_1[iT_1] \rightsquigarrow t'_2} \quad (\text{PE-INLINETAPP}) \\
\frac{\text{inline } v \rightsquigarrow \text{inline } v}{t \rightsquigarrow \text{inline } t'} \quad (\text{PE-INLINEVALUE}) \\
\frac{t \rightsquigarrow \text{inline } t'}{\text{dynamic } t \rightsquigarrow t'} \quad (\text{PE-DYNAMIC})
\end{array}$$

Fig. 7. Partial evaluation $t \rightsquigarrow t'$

2.6 Evaluation

Once partial evaluation is complete we strip all binding-time annotations on types and convert inline terms into corresponding dynamic ones. After that we can use standard $F_{<}$ evaluation rules augmented with lazy records semantics (E-SEL).

$$\begin{array}{c}
\frac{v \Downarrow v}{t_1 \Downarrow (x : T) \Rightarrow t \quad t_2 \Downarrow v \quad [x \mapsto v]t \Downarrow v'} \quad (\text{E-VALUE}) \\
\frac{t_1(t_2) \Downarrow v'}{t_1 \Downarrow [X <: T_2] \Rightarrow t_2 \quad [X \mapsto T_1]t_2 \Downarrow v} \quad (\text{E-APP}) \\
\frac{t_1[T_1] \Downarrow v}{t \Downarrow \{x = t_x, \overline{y = t_y}\} \quad t_x \Downarrow v} \quad (\text{E-TAPP}) \\
\frac{t \Downarrow \{x = t_x, \overline{y = t_y}\} \quad t_x \Downarrow v}{t.x \Downarrow v} \quad (\text{E-SEL})
\end{array}$$

Fig. 8. Evaluation $t \Downarrow v$

2.7 Conjectures

1. Progress and preservation of partial evaluation.
2. Progress and preservation of evaluation.

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

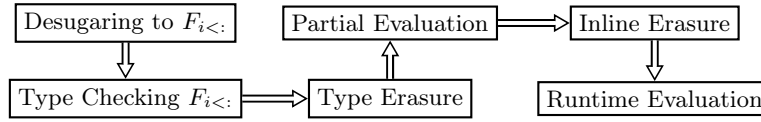


Fig. 9. Compilation pipeline.

The $F_{i<}$ calculus §2 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 $F_{i<}$ (§3.1), show how to provide compile time views of classes and *methods*§??, and formalize convenient implicit conversions for the calculus §3.3.

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.

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

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

Fig. 10. Desugaring of classes into $F_{i<}$.

3.2 Compile-Time View of the Terms

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

Fig. 11. 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.

3.3 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.

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
```

4 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 (§4.2), then we demonstrate how variable argument functions can be desugared into the core functionality (§4.3). Finally, we demonstrate how the abstraction overhead of the `dot` function and all associated type-classes can be removed (§4.5).

4.1 Inlining Expressed Through Partial-Evaluation

4.2 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]**

4.3 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 12.

```
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. 12. 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`.

4.4 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.

4.5 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.

5 The Partial Evaluator for Scala

```

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

[TODO: cite github] We have implemented a prototype partial evaluator §?? and the desugaring §3 for the Scala language. 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 bidding-time of terms.

To the user, the partial evaluator exposes a minimal interface (Figure 5) with `inline` and `ct` annotations and the `ct` function.

```

package object scalainline {
  final class ct extends StaticAnnotation

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

  final class inline extends StaticAnnotation
}

```

Fig. 14. Interface of the Scala partial evaluator.

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

In Table 2, `Int@ct` is a non-polymorphic type and therefore according to the translation to the compile-time view (11) all arguments of all methods will be compile-time. On the other hand, `Vector[Int]@ct` will have all arguments

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>

of all methods transformed except the generic ones. Function `f` passed to `map` accepts a dynamic value as input.

& All operations on `Int` executed at compile-time.

& `\code{map}` executed at compile-time, over dynamic values.

& `Length` is executed at compile-time; result is compile-time.

& `\code{map}` executed at compile-time over compile-time values. The result can still be both d

& `Map`

Method `ct` is used at the term level for promoting Scala objects and functions into their compile-time views.

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>ct((x: Int) => x)</code>	<code>: (Int => Int)@ct</code>
<code>ct((x: Int@ct) => x)</code>	<code>: (Int@ct => 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>

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

5.1 Interaction with the Scala Language

6 Evaluation

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 Wölß, 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.