

Automating Ad hoc Data Representation Transformations

Abstract

To maximize run-time performance, programmers often specialize their code by hand, replacing library collections and containers by custom objects in which the data are restructured for efficient access. However, this is a tedious and error-prone process that makes it hard to test, maintain and evolve the source code.

We present an automated and composable approach that allows programmers to safely change the data representation in delimited scopes containing anything from expressions to entire class definitions. The transformation itself is programmer-defined and covers a wide range of use cases.

Our technique leverages the type system in order to infer where the data representation needs to be converted, while offering a strong correctness guarantee on the interaction with object-oriented language features, such as dynamic dispatch, inheritance and generics.

We have embedded this technique in a Scala compiler plugin and used it in four very different transformations, ranging from improving the data layout and encoding, to retrofitting specialization and value class status, to collection deforestation. As a result, we obtain large speedups: between 1.9x and 14.5x.

1. Introduction

An object encapsulates code and data and exposes an interface. Modern language facilities, such as extension methods, type classes and implicit conversions allow programmers to evolve the object interface in an ad hoc way, by adding new methods and operators. For example, in Scala, we can use an implicit conversion to add the multiplication operator to pairs of integers, with the semantics of complex number multiplication:

```
1 scala> (0, 1) * (0, 1)
2 res0: (Int, Int) = (-1, 0)
```

Unlike evolving the interface, no general mechanism in modern languages is capable of evolving an object’s encapsulated data. The encapsulated data format is assumed to be fixed, allowing the compiled code to contain hard references to data, encoded according to a convention known as the *object layout*. For instance, methods encapsulated by the generic pair class, such as `swap` and `toString`, rely on the existence of two generic fields, erased to `Object`. This leads to inefficient storage in the running example, as the integers need to be boxed, producing as many as 3 heap objects for each “complex number”: the two boxed integers

and the pair container. What if, for a part of our program, instead of the pair, we concatenated the two 32-bit integers into a 64-bit long integer, that would represent the “complex number”? We could pass complex numbers by value, completely sidestepping the need to allocate memory and to garbage-collect it later. Additionally, what if we could also add functionality, such as arithmetic operations, to our ad hoc complex numbers, all without any heap allocation overhead? Finally, what parts of such a transformation could be automated?

Object layout transformations are common in language virtual machines, such as V8 and Truffle. These virtual machines profile values at run-time and make optimistic assumptions about the shape of objects. This allows them to automatically optimize the object layout in the heap, at the cost of recompiling of all the code that references the old object layout. If, later in the execution, the assumptions prove too optimistic, the virtual machine needs to revert to the more general (and less efficient) object layout, again recompiling all the code that contains hard references to the optimized layout. As expected, this comes with significant overheads. Thus, runtime decisions to change the low-level layout are expensive (due to recompilation) and have a global nature, affecting all code that assumes a certain layout.

Since transforming the object layout at run-time is expensive, a natural question to ask is whether we can leverage the statically-typed nature of a programming language to optimize the object layout during compilation? The answer is yes. Transformations such as “class specialization” and “value class inlining” transform the object layout in order to avoid the creation of heap objects. However, both of these transformations take a global approach: when a class is marked as specialized or as a value class (and assuming it satisfies the semantic restrictions) it is transformed at its definition site. Later on, this allows all references to the class, even in separately compiled sources, to be optimized. On the other hand, if a class is not marked at its definition site, retrofitting specialization or the value class status is impossible, as it would break many non-orthogonal language features, such as dynamic dispatch, inheritance and generics.

Therefore, although transformations in statically typed languages can optimize the object layout, they do not meet the ad hoc criterion: they cannot be retrofitted later, and they have a global, all-or-nothing nature. For instance, in Scala, the generic pair class is specialized but not marked as a value class. As a result, the representation is not fully optimized, still requiring a heap object per pair. Even worse, specializa-

tion and value class inlining are mutually exclusive, making it impossible to optimally represent our “complex numbers” as values even if we had complete control over the Scala library. Furthermore, a change in the data representation may be applicable for specific parts of the client code, but might not make sense to apply globally.

In our “complex numbers” abstraction, we only use a fraction of the flexibility provided by the library tuples, and yet we have to give up all the code optimality. Even worse, for our limited domain, we are aware of a better representation, but the only solution is to transform the code by hand, essentially having to choose between an obfuscated or a slow version of the code. What is missing is a largely automated and safe transformation that allows us to use our domain-specific knowledge to mark a scope where the “complex numbers” can use the alternative object layout, effectively specializing that part of our program.

In this paper we present such an automated transformation that allows programmers to safely change the data representation in limited, well-defined scopes that can include anything from expressions to method and class definitions. The transformation maintains strong correctness guarantees in terms of non-orthogonal language features, such as dynamic dispatch, inheritance and generics across separate compilations. To gain the most benefit, our approach uses the programmers’ domain-specific knowledge of the transformed scope, allowing them to specify the exact alternative representation and the operations it should expose, while automating all the tedium involved in safely transforming the code and maintaining the outside interface.

In this way, the programmer is responsible for correctly stating (a) what the data representation transformation is and (b) to which program scope it applies. Our approach is then responsible for (1) automatically deciding when to apply the transformation and when to revert it, in order to ensure correct interchange between representations, (2) enriching the transformation with automatically generated bridge code that ensures correctness relative to overriding and dynamic dispatch and (3) persisting the necessary metadata to allow transformed program scopes in different source files and compilation runs to communicate using the optimized representation—a property we refer to as *composability* in the following sections. Thus, our approach adheres to the design principle of separating the reusable, general and provably correct *mechanism* from the programmer-defined *policy*, which may contain incorrect decisions [26].

Our main contributions are:

- Introducing the ad hoc data representation problem, which, to the best of our knowledge, has not been addressed at all in the literature (§2);
- Presenting the extensions that allow global data representation transformations (§3) to be used as scoped programmer-driven transformations (§4);
- Implementing the approach presented as a Scala compiler plugin [2] that allows programmers to express custom transformations (§5) and benchmarking the plugin on a

broad spectrum of transformations, ranging from improving the data layout and encoding, to retrofitting specialization and value class status, and to collection deforestation [44]. These transformations produced speedups between 1.9 and 14.5x on user programs (§6).

2. Motivation

This section presents the full motivating example featuring the complex number transformation, which we use throughout the paper. It then shows how the data representation transformation is triggered and introduces the terminology. Finally, it shows a naive transformation, hinting at the difficulties lying ahead.

2.1 Motivating Example

In the introduction, we focused on adding complex number semantics to pairs of integers. Complex numbers with integers as both their real and imaginary parts are known as Gaussian integers, and are a countable subset of all complex numbers. The operations defined on Gaussian integers are similar to complex number operations, with one exception: to satisfy the abelian closure property, division is not precise, but instead rounds the result to the nearest Gaussian integer, with both the real and imaginary axes containing integers. This is similar to integer division, which also rounds the result, so that, for example, $5/2$ produces value 2.

An interesting property of Gaussian integers is that we can define the “divides” relation and the greatest common divisor (GCD) between any two Gaussian integers. Furthermore, computing the GCD is similar to Euclid’s algorithm for integer numbers:

```
1 def gcd(n1: (Int, Int), n2: (Int, Int)): (Int, Int) = {
2   val remainder = n1 % n2
3   if (remainder.norm == 0) n2 else gcd(n2, remainder)
4 }
```

Unfortunately, as our algorithm recursively computes the result, it creates linearly many pairs of integers, allocating them on the heap. If we run this algorithm with no optimizations, computing the GCD takes around 3 microseconds (on the same setup as used for our full experiments in §6):

```
1 scala> timed(() => gcd((544, 185), (131, 181)))
2 The operation took 3.05 us (based on 10000 executions).
3 The result was (10, 3).
```

However, if we encode the Gaussian integers into 64-bit long integers, we improve the time by a factor of 13x:

```
1 scala> timed(() => gcd((544, 185), (131, 181)))
2 The operation took 0.23 us (based on 10000 executions).
3 The result was (10, 3).
```

This makes a data representation transformation highly desirable. Still, making the programmer transform the code by hand is tedious and error-prone, so a natural question to ask is whether the transformation could be automated.

2.2 Automating the Transformation

In order to reap the benefits of using the improved representation without manually transforming the code, we

present the Ad hoc Data Representation (ADR) Transformation technique, which can be triggered by the `adrt` marker method. This method accepts two parameters: the first parameter is a *transformation description object* and the second is a block of code that forms the transformation scope. This scope can contain anything from expressions all the way to method, class, trait and object definitions:

```
1 adrt(IntPairComplexToLongComplex) {
2   def gcd(n1: (Int, Int), n2: (Int, Int)): (Int, Int) = {
3     val remainder = n1 % n2
4     if (remainder.norm == 0) n2 else gcd(n2, remainder)
5   }
6 }
```

To maintain a consistent naming throughout the paper, we will use the name *high-level type* to designate `(Int, Int)`, which corresponds to the original type in the code. This high-level type can be encoded as its *representation type*, `Long`. The high-level type, its representation, and the procedures for encoding and decoding are all stored in the transformation description object, in our case `IntPairComplexToLongComplex`.

2.3 A Naive Transformation

Despite its simple interface, the Ad hoc Data Representation Transformation mechanism is by no means simple. For example, knowing that the `adrt` marker was instructed to transform `(Int, Int)` to `Long`, a naive result could be:

```
1 def gcd(n1: Long, n2: Long): Long = {
2   val remainder = n1 % n2
3   if (remainder.norm == 0) n2 else gcd(n2, remainder)
4 }
```

There are many questions one could ask about this naive translation. For example, how does the compiler know which parameters and values to transform to the long integer representation (§4.1)? How and when to encode and decode values, and what to do about values that are visible outside the scope (§4.2)? Even worse, what if parts of the code are compiled separately, in a different compiler run (§4.3)?

Going into the semantics of the program, we can ask if the `%` (modulo) operator maintains the semantics of complex numbers when used for long integers. Also, is `norm` defined for long integers? Unfortunately, the response to both questions is negative. Therefore, the transformation needs to preserve semantics, which is not trivial (§4.4).

We could also ask what would happen if `gcd` was overriding another method, as in the code fragment below. Would the new signature still override it? The answer is no, so the naive translation would break the object model (§4.5).

```
1 trait WithGCD[T] {
2   def gcd(n1: T, n2: T): T
3 }
4 class Complex extends WithGCD[(Int, Int)] {
5   // expected: gcd(n1: (Int, Int), n2: (Int, Int)) ...
6   // found: gcd(n1: Long, n2: Long): Long
7   // (which does not implement gcd in trait WithGCD)
8   def gcd(n1: Long, n2: Long): Long = ...
9 }
```

Our approach, the Ad hoc Data Representation Transformation, addresses these questions.

3. Data Representation Transformations

As necessary background for our approach, we review data representation transformations in general and, in particular, the Late Data Layout (LDL) transformation mechanism [42], which we later extend to our ad hoc data representation transformation (§4). The LDL mechanism is neither programmer-driven, since the data representation has to be known a priori and encoded in the transformation, nor directly applicable to limited scopes inside a program.

Data can usually be represented in several ways, some more efficient and others more flexible. For example, integer numbers can use either the primitive (unboxed) value encoding, which is more efficient, or the object-based (boxed) encoding, which is more flexible. The boxed representation allows integers to act as the receivers of dynamically dispatched method calls, to be assigned to supertypes, such as `Number` or `Object` and to instantiate erased generics. However, the extra flexibility comes at a price: boxed integers are allocated on the heap so they need to be garbage-collected later and all their operations incur an indirection overhead. This leads to a tension between the two representations.

From a language perspective, there are two approaches to exposing the multiple representations of a type: either have a different type for each representation, as Java does, or fully hide the difference and present a single language-level type, as ML, Haskell and Scala do. Either way, the final low-level bytecode or assembly code needs to handle the two representations separately, since they correspond to very different entities: references and values.

Exposing a single high-level type in the language is more popular among programmers for its simplicity, but it places more responsibility on the compiler, which has to perform two additional steps: first, it needs to choose the data representation of each value; and second, it needs to introduce coercions that switch between representations where necessary. For example, since only boxed integers can instantiate generics, any unboxed integer going into a generic container, such as a list, needs to be *coerced* to the boxed representation. This work is done in the compiler pipeline, in so-called data representation transformations.

The Late Data Layout mechanism, presented next, is a powerful data representation transformation facility for Scala. It has three properties that make it well-suited to be a substrate for our Ad hoc Data Representation Transformation: selectivity, optimality and consistency.

3.1 Late Data Layout

The Late Data Layout (LDL) mechanism [42] is the underlying transformation used in Scala to implement multi-parameter value class inlining and to specialize classes using the miniboxed encoding [41]. It is a flexible and reliable mechanism, tested on many thousands of lines of code.

Using LDL, a language can expose high-level types (called high-level concepts in the LDL terminology), such as the integer type `Int` exposed by Scala, which can represent either a boxed or unboxed value in the low-level bytecode. In

the following example, we have values of types `Int` and `Any`. `Any` is the top of the Scala type system, and thus a supertype of `Int`:

```
1 val i: Int = 1
2 val j: Int = i
3 val k: Any = j
```

Since Scala compiles down to Java bytecode, during compilation, the LDL-based primitive unboxing transformation bridges the gap between the high-level `Int` concept and its two representations: the unboxed `int` and the boxed `java.lang.Integer` representation. Along the way, it introduces the necessary coercions between these two representations. For example, the code above is translated to:

```
1 val i: int = 1
2 val j: int = i
3 val k: Any = Integer.valueOf(j)
```

The LDL mechanism transforms the data representation in three phases: **INJECT**, **COERCE** and **COMMIT**. Each of the phases is responsible for a property of the transformation: **INJECT** makes LDL *selective*, **COERCE** makes it *optimal* and **COMMIT** makes it *consistent*. In our examples, we show the equivalent source code for the program abstract syntax trees (ASTs) after each of these phases.

The INJECT phase is responsible for marking each value with its desired representation. In the case of primitive integer unboxing, the annotation is `@unboxed`, and it signals that a value should be stored in the unboxed `int` representation. As an optimization, instead of adding a `@boxed` annotation for the corresponding cases, values that are not marked are automatically considered boxed. Following the **INJECT** phase, the previous example will be transformed to:

```
1 val i: @unboxed Int = 1 // Int can be unboxed
2 val j: @unboxed Int = i // Int can be unboxed
3 val k: Any = j          // Any cannot be unboxed
```

The **INJECT** phase gives LDL a selective nature, allowing it to mark each individual value with its representation. For example, it would have been equally correct if the marking rules decided that `j` should be boxed, in which case its type would not have been marked. One of the properties of the LDL transformation is that boxed and unboxed values are compatible in the **INJECT** phase, so there are no coercions.

The COERCE phase, as its name suggests, introduces coercions. This is done by changing the annotation semantics: annotated types become incompatible among themselves and with their un-annotated counterparts. This change in the annotation semantics corresponds to introducing the different representations: each annotation corresponds to a representation, and representations are not compatible with each other. With this change, an assignment from one representation to another will lead to mismatching types. Therefore, by re-type-checking the tree, the **COERCE** phase can detect representation mismatches and can patch them using coercions. In the example, the last line contains such a mismatch:

```
1 val i: @unboxed Int = 1 // expected/found: @unboxed
2 val j: @unboxed Int = i // expected/found: @unboxed
3 val k: Any = box(j)    // mismatch => box
```

The **COERCE** phase establishes the optimality property of the LDL transformation. The definition of optimality is quite involved, but we can easily show it using an example. Consider the following two integer definitions:

```
1 val c: Boolean = ...
2 val l1: @unboxed Int = if (c) i else j
3 val l2: @unboxed Int = unbox(if (c) box(i) else box(j))
```

It is clear that the two definitions will always produce the same result. Yet, the first one is markedly better: it does not execute any coercions, compared to second definition, which executes two coercions regardless of the value of `c`. These subtle sub-optimality can slow down program execution, increase the heap footprint and the bytecode size. This is why the **COERCE** phase needs to be optimal. The LDL paper [42] makes the following intuition-based conjecture: “in any given terminating execution trace through the transformed program, the number of coercions executed is minimal, for given sets of annotations introduced by the **INJECT** phase and transformations performed in the **COMMIT** phase”. An initial formalization and proof is sketched in [40].

From our perspective, optimality means that once representations are chosen and annotated, only the necessary coercions will be introduced during the **COERCE** phase.

The COMMIT phase is responsible for introducing the actual representations. In the case of primitive unboxing, `@unboxed Int` is replaced by `int`, and `Int`, which is considered boxed, is replaced by `java.lang.Integer`. The `box` and `unbox` coercions are also replaced by the creation of objects and, respectively, by the extraction of the unboxed value:

```
1 val i: int = 1
2 val j: int = i
3 val k: Any = Integer.valueOf(j)
```

The **COMMIT** phase is responsible for the consistency of the transformation. Since the program abstract syntax tree (AST) has been checked by the type-system extended with representation semantics, the **COMMIT** phase is guaranteed to correctly handle the value representations and to correctly coerce between them. This allows the **COMMIT** phase to be a very simple, syntax-based, transformation over the program abstract syntax tree (AST).

3.2 Support For Object-Oriented Programming

The LDL mechanism targets object-oriented programming languages, which pose unique challenges for data representation transformations. This section will describe the additional rules necessary in LDL to handle object-orientation.

Object-oriented Patterns. Aside from introducing coercions, data representation transformations must handle object-oriented patterns, such as method calls and subtyping. Not all representations can be used with these patterns. For example, it is not possible to call the `toString` method on the unboxed `int` representation:

```
1 val a: @unboxed Int = 1
2 println(a.toString)
```

To handle dynamically dispatched method calls, LDL has a built-in rule: when a value acts as a method call receiver, it is coerced to the boxed representation, which, in this case, corresponds to the non-annotated representation. In our example, the `@unboxed Int` value is boxed during the COERCE phase, so it can act as the receiver of the `toString` method:

```
1 val a: @unboxed Int = 1
2 println(box(a).toString)
```

To improve performance, the LDL mechanism also supports bypass methods, also known as extension methods in the literature. For example, if a static `bypass_toString` method is available for the unboxed `int` representation, there is no need to convert it before the method call:

```
1 val a: @unboxed Int = 1
2 println(bypass_toString(a))
```

Subtyping is handled in a similar fashion, by requiring the boxed representation, which can be assigned to supertypes.

Support for Generics. The Late Data Layout mechanism is agnostic to generics. This means that, depending on the transformation semantics and the implementation of generics, the mechanism can inject annotations in the type arguments or not. For example, if generics are erased, a list of integers will have type `List[Int]`, since values need to be boxed. If generics are unboxed and reified, the list type will be `List[@unboxed Int]`. In the LDL paper [42], the authors show examples of both cases: when annotations are propagated inside generics and when they are not. The LDL mechanism adapts seamlessly to either case.

Having seen the Late Data Layout mechanism at work for unboxing primitive types, we can now look at how it can be extended to handle ad hoc programmer-driven data representation transformations.

4. Ad hoc Data Representation Transformation

The Ad hoc Data Representation (ADR) transformation adds two new elements to existing data representation transformations: (1) it enables custom, programmer-defined alternative representations and (2) it allows the transformation to take place in limited scopes, ranging from expressions all the way to method and class definitions. This allows programmers to use locally correct transformations that may be incorrect for code outside their given scope.

Section 2.2 showed how the ADR transformation is triggered by the `adrt` marker. The running example is reproduced below for quick reference:

```
1 adrt(IntPairComplexToLongComplex) {
2   def gcd(n1: (Int, Int), n2: (Int, Int)): (Int, Int) = {
3     val remainder = n1 % n2
4     if (remainder.norm == 0) n2 else gcd(n2, remainder)
5   }
6 }
```

The following sections take a step by step approach to explaining how our technique allows programmers to define transformations and to use them in localized program scopes, improving the performance of their programs in an automated and safe fashion.

4.1 Transformation Description Objects

The first step in performing an `adrt` transformation is defining the transformation description object. This object is required to extend a marker interface and to define the transformation through the `toRepr` and `toHigh` coercions:

```
1 object IntPairComplexToLongComplex
2   extends TransformationDescription {
3   // coercions:
4   def toRepr(high: (Int, Int)): Long = ...
5   def toHigh(repr: Long): (Int, Int) = ...
6   // bypass methods:
7   ...
8 }
```

The coercions serve a double purpose: (1) the signatures match the high-level type, in this case `(Int, Int)` and indicate the corresponding representation type, `Long` and vice-versa and (2) the implementations are called in the transformed scope to encode and decode values as necessary.

Bypass Methods. The description object can optionally include bypass methods, which correspond to the methods exposed by the high-level type, but instead operate on encoded values in the representation type. Bypass methods allow the transformation to avoid coercing receivers to the high-level type by rewriting dynamically dispatched calls to their corresponding statically-resolved bypass method calls, as shown in section §3.2. Method call rewriting in `adrt` scopes is explained later, in section §4.4.

Generic Transformations. In our example, both the high-level and representation types are monomorphic (i.e. not generic). Still, in some cases, the ADR transformation is used to target collections regardless of the type of their elements. We analyzed multiple approaches to allowing genericity in the transformation description object and converged on allowing the coercions to be generic themselves. This approach has the merit of being concise and extending naturally to any arity:

```
1 def toRepr[T](high: List[T]): LazyList[T] = ...
2 def toHigh[T](repr: LazyList[T]): List[T] = ...
```

Since the coercion signatures “match” the high-level type and return the corresponding representation type, a value of type `List[Int]` will be matched by the `adrt` transformation and subsequently encoded as a `LazyList[Int]`. This allows the `adrt` scopes to transform collections, containers and function representations. The benchmarks section (§6) shows two examples of generic transformations.

Target Semantics. It is worth noting that coercions defined in transformation objects must maintain the semantics of the high-level type. In particular, semantics such as mutability and referential identity must be preserved if the program relies on them. For example, correctly handling referential

identity requires the coercions to return the exact same object (the exact same reference) when interleaved:

```
1 assert(toHigh(toRepr(x)) eq x) // referential equality
```

These semantics reduce the benefit of the `adrt` transformation by imposing restrictions on the coercions. However, in most use cases, the targets, such as library collections and containers, have value semantics: they are immutable, final and only use structural equality. Such high-level types can be targeted at will, since they can be reconstructed at any time without the program observing it.

Once the transformation description object is defined, it can be used in `adrt` scopes to optimize the user program.

4.2 Transformation Scopes and Composability

Existing data representation transformations, such as value class inlining and specialization, have fixed semantics and occur in a sequence. Instead, the ADR transformation handles all transformation scopes in the source code concurrently, each with its own high-level target, representation type and coercions. This is a challenge, as handling the interactions between these concurrent scopes, some of which may even be nested, demands a highly disciplined treatment.

The key to handling all concurrent scopes correctly is shifting focus from the scopes themselves to the values they define. Since we are using the underlying LDL mechanism, we can track the encoding of each value in its type, using annotations. To keep track of the different transformations introduced by different scopes, we extend the LDL annotation to reference the description object, essentially carrying the entire transformation semantic with each individual value. We then leverage the type system and the signature persistence facilities to correctly transform all values, essentially allowing scopes to safely and efficiently pass values among themselves, using the representation type—a property we refer to as composability.

We look at four instances of composability:

- allowing different scopes to communicate, despite using different representation types (high-level types coincide);
- isolating high-level types, barring unsound value leaks through the representation type;
- handling conflicting transformation description objects;
- passing values between high-level types in the encoded (representation) format;

Although the four examples cover the most interesting corner cases of the transformation, the interested reader may consult the “Scope Nesting” page on the project wiki [3], which describes all cases of scope overlapping, collaboration and nesting. Furthermore, the scope composition is tested with each commit, as part of the project’s test suite.

A high-level type can have different representations in different scopes. This follows from the scoped nature of the ADR transformation, which allows programmers to use the most efficient data representation for each task. But it raises the question of whether values can be safely passed across scopes that use different representations:

```
1 adrt(IntPairToLong) { var x = (3, 5) }
2 adrt(IntPairToDouble) { val y = (2, 6); x = y }
```

At a high level, the code is correct: the variable `x` is set to the value of `y`, both of them having high-level type `(Int, Int)`. However, being in different scopes, these two values will be encoded differently, `x` as a long integer and `y` as a double-precision floating point number. In this situation, how will the assignment `x = y` be translated? Let us look at the transformation step by step.

After parsing, the scope is inlined and the program is type-checked against the high-level types. Aside from checking the high-level types, the type checker also resolves implicits and infers all missing type annotations. During type checking, the description objects are stored as invisible abstract syntax tree attachments (described in §5):

```
1 var x: (Int, Int) = (3, 5) /* att: IntPairToLong */
2 val y: (Int, Int) = (2, 6) /* att: IntPairToDouble */
3 x = y
```

Then, during the INJECT phase, each value or method definition that matches the description object’s high-level type is annotated with the `@repr` annotation, parameterized on the transformation description object:

```
1 var x: @repr(IntPairToLong) (Int, Int) = (3, 5)
2 val y: @repr(IntPairToDouble) (Int, Int) = (2, 6)
3 x = y
```

The `@repr` annotation is only attached if the value’s type matches the high-level type in the description object. Therefore, programmers are free to define values of any type in the scope, but only those values whose type matches the transformation description object’s target will be annotated.

Based on the annotated types, the COERCE phase notices the mismatching transformation description objects in the last line: the left-hand side is on its way to be converted to a long integer (based on the description object `IntPairToLong`) while the right-hand side will become a floating point expression (based on the description object `IntPairToDouble`). However, both description objects have the same high-level type, the integer pair. Therefore, the high-level type is used as a middle ground to transform between the two representation types:

```
1 var x: @repr(IntPairToLong) (Int, Int) =
  toRepr(IntPairToLong, (3, 5))
2 val y: @repr(IntPairToDouble) (Int, Int) =
  toRepr(IntPairToDouble, (2, 6))
3 x = toRepr(IntPairToLong, toHigh(IntPairToDouble, y))
```

Finally, the COMMIT phase transforms the example to:

```
1 var x: Long = IntPairToLong.toRepr((3, 5))
2 val y: Double = IntPairToDouble.toRepr((2, 6))
3 x = IntPairToLong.toRepr(IntPairToDouble.toHigh(y))
```

In the end, the value `x` is converted from a double to a pair of integers, which is subsequently converted to a long integer. This shows the disciplined way in which different `adrt` scopes compose, allowing values to flow across different representations, from one scope to another. As in the LDL transformation, the mechanism aims to employ a minimal number of conversions.

Different transformation scopes can be safely nested and the high-level types are correctly isolated:

```
1 adrt (FloatPairAsLong) {
2   adrt (IntPairAsLong) {
3     val x: (Float, Float) = (1f, 0f)
4     var y: (Int, Int) = (0, 1)
5     // y = x
6     // y = 123.toLong
7   }
8 }
```

Values of the high-level types in the inner scope are independently annotated and are transformed accordingly. Since both the integer and the float pairs are encoded as long integers, a natural question to ask is whether values can leak between the two high-level types, for example, by uncommenting the last two lines of the inner scope. This would open the door to incorrectly interpreting an encoded value as a different high-level type, introducing unsoundness.

The answer is no: the code is first type-checked against the high-level types even before the INJECT transformation has a chance to annotate it. This prohibits direct transfers between the high-level types and their representations. Thus, the unsound assignments will be rejected, informing the programmer that the types do not match. This is a non-obvious benefit of using the ADR transformation instead of manually refactoring the code and using implicit conversions, which would allow such unsound assignments.

Handling conflicting nested transformation description objects is another important property of composition:

```
1 adrt (PairAsMyPair) {
2   adrt (IntPairAsLong) {
3     val x: (Int, Int) = (2, 3)
4   }
5   println(x.toString)
6 }
```

In the code above, the type of `x` matches both transformation description objects, so it could be transformed to both representation types `MyPair[Int, Int]` and `Long`. However, during the INJECT phase, if a value is matched by several nested `adrt` scopes, this can be reported to the programmer either as an error or, depending on the implementation, as a warning, followed by choosing one of the transformation description objects for the value (current solution):

```
1 console:9: warning: Several adrt scopes can be applied
   to value x. Picking the innermost one: IntPairAsLong
2 val x: (Int, Int) = (2, 3)
3   ^
```

Furthermore, since the INJECT phase annotates value `x` with the chosen transformation, there will be no confusion in the next line, where `x` has to be converted back to the high-level type to receive the `toString` method call, despite the fact that the `adrt` scope surrounding the instruction uses a different transformation description object.

Prohibiting access to the representation type inside the transformation scope is limiting. For example, a performance-conscious programmer might want to transform the high-level integer pair into a floating-point pair without allocating heap objects. Since the programmer does not have

direct access to the representation, it looks like the only solution is to decode the integer pair into a heap object, convert it to a floating-point pair and encode it back to the long integer.

There is a better solution. As we will later see, the programmer can use bypass methods to “serialize” the integer pair into a long integer and “de-serialize” it into a floating-point pair. Yet, this requires a principled change in the transformation description object. This is the price to pay for a safe and automated representation transformation.

To recap: focusing on individual values and storing the transformation semantics in the annotated type allows us to correctly handle values flowing across scopes, a property we call scope composition. Although we focused on values, method parameters and return types are annotated in exactly the same way. The next part extends scope composition across separate compilation.

4.3 Separate Compilation

Annotating the high-level type with the transformation semantics allows different `adrt` scopes to seamlessly pass encoded values. To reason about composing scopes across different compilation runs, let us assume we have already compiled the `gcd` method in the motivating example:

```
1 adrt (IntPairComplexToLongComplex) {
2   def gcd(n1: (Int, Int), n2: (Int, Int)): (Int, Int) = ..
3 }
```

After the INJECT phase, the signature for method `gcd` is:

```
1 def gcd(
2   n1: @repr(IntPairComplexToLongComplex) (Int, Int),
3   n2: @repr(IntPairComplexToLongComplex) (Int, Int)
4 ): @repr(IntPairComplexToLongComplex) (Int, Int) =
   ...
```

And, after the COMMIT phase executed, the bytecode signature for method `gcd` is:

```
1 def gcd(n1: long, n2: long): long = ...
```

When compiling source code that refers to existing low-level code, such as object code or bytecode compiled in a previous run, compilers need to load the signature of each symbol. For C and C++ this is done by parsing header files while for Java and Scala, it is done by reading the source-level signature from the bytecode metadata. However, not being aware of the ADR transformation of method `gcd`, a separate compilation could assume it accepts two pairs of integers as input. Yet, in the bytecode, the `gcd` method accepts longs and is not able to handle pairs of integers.

The simplest solution is to create two versions for each transformed method: the transformed method itself and a bridge, which corresponds to the high-level signature. The bridge method would accept pairs of integers and encode them as longs before calling the transformed version of the `gcd` method. It would also decode the result of `gcd` back to a pair of integers. This approach allows calling `gcd` from separately compiled files without being aware of the transformation. Still, we can do better.

Persisting transformation annotations. Let us assume we want to call the `gcd` method from a scope transformed using

the same transformation description object as we used when compiling `gcd`, but in a different compilation run:

```
1 adrt (IntPairComplexToLongComplex) {
2   val n1: (Int, Int) = ...
3   val n2: (Int, Int) = ...
4   val res: (Int, Int) = gcd(n1, n2)
5 }
```

In this case, would it make sense to call the bridge method? The values `n1` and `n2` are already encoded, so they would have to be decoded before calling the bridge method, which would then encode them back. This is suboptimal. Instead, what we want is to let the `adrt` scopes become part of the high-level signature, but without making the transformation a first-class language feature. To do this, instead of persisting the scope, we persist the injected annotations, including the reference to the transformation description object. These become part of the signature of `gcd`:

```
1 // loaded signature (description object abbreviated):
2 def gcd(n1: @repr(.) (Int, Int), n2: @repr(.) (Int,
   Int)):@repr(.) (Int, Int)
```

The annotations are loaded just before the INJECT phase, which transforms our code to:

```
1 val n1: @repr(.) (Int, Int) = ...
2 val n2: @repr(.) (Int, Int) = ...
3 val res: @repr(.) (Int, Int) = gcd(n1, n2)
```

With the complete signature for `gcd`, the COERCE phase does not introduce any coercions, since the arguments of method `gcd` are going to be transformed in the same way as the method parameters were transformed in a previous compilation run. This allows `adrt` scopes to seamlessly compose even across separate compilation. After the COMMIT phase, the scope is compiled to:

```
1 val n1: Long = ...
2 val n2: Long = ...
3 val res: Long = gcd(n1, n2) // no coercions!!!
```

Making bridge methods redundant. Persisting transformation information in the high-level signatures allows us to skip creating bridges. For example, calling the `gcd` method outside the `adrt` scope is still possible:

```
1 val res: (Int, Int) = gcd((55, 2), (17, 13))
```

Since the signature for method `gcd` references the transformation description object, the COERCE phase knows exactly which coercions are necessary:

```
1 val res: (Int, Int) = toHigh(...,
2   gcd(toRepr(..., (55, 2)), toRepr(..., (17, 13))))
```

Generally, persisting references to the description objects in each value's signature allows scope composition across separate compilation runs.

4.4 Optimizing Method Invocations

When choosing a generic container, such as a pair or a list, programmers are usually motivated by the very flexible interface, which allows them to quickly achieve their goal by invoking the container's many convenience methods. The presentation so far focused on optimizing the data represen-

tation, but to obtain peak performance, the method invocations need to be transformed as well:

```
1 adrt (IntPairComplexToLongComplex) {
2   val n = (0, 1)
3   println(n.toString)
4 }
```

When handling method calls on an encoded receiver, the default LDL behavior is very conservative: it decodes the value back to its high-level type, which exposes the original method and generates a dynamically-dispatched call (§3.2):

```
1 val n: Long = ...
2 println(IntPairComplexToLongComplex.toHigh(n).toString)
```

The price to pay is decoding the value into the high-level type, which usually leads to heap allocations and can introduce overheads. If a corresponding bypass method is available, the LDL transformation can use it:

```
1 val n: Long = ...
2 println(IntPairComplexToLongComplex.bypass_toString(n))
```

The bypass method can operate directly on the encoded version of the integer pair, avoiding a heap allocation. In practice, when the receiver of a method call is annotated, our modified LDL transformation looks up the `bypass_toString` method in the transformation description object, and, if none is found, warns the programmer and proceeds with decoding the receiver and generating the dynamically-dispatched call.

Methods added via implicit conversions and other enrichment techniques, such as extension methods or type classes, add another layer or complexity, only handled in the ADR transformation. For example, we can see the multiplication operator `*`, added via an implicit conversion (we will further analyze the interaction with implicit conversions in §4.5):

```
1 adrt (IntPairComplexToLongComplex) {
2   val n1 = (0, 1)
3   val n2 = n1 * n1
4 }
```

Type-checking the program introduces an explicit call to the implicit conversion that adds the correct `*` operator:

```
1 val n1: (Int, Int) = (0, 1)
2 val n2: (Int, Int) = intPairAsComplex(n1) * n1
```

This is a costly pattern, requiring `n1` to be decoded into a pair and passed to the `intPairAsComplex` method, which itself creates a wrapper object that exposes the `*` operator. To optimize this pattern, the ADR transformation looks for a bypass method in the transformation description object that corresponds to a mangled name combining the implicit method name and the operator. For simplicity, if we assume the name is `implicit_*` and the bypass exists, we get:

```
1 val n1: Long = IntPairComplexToLongComplex.toRepr(0, 1)
2 val n2: Long = implicit_*(n1, n1)
```

This allows the call to the `*` operator to go through without any heap object creation, which significantly improves performance and heap footprint. We next see the implementation of bypass methods in detail.

Bypass methods. Both normal and implicit bypass methods need to correspond to the method they are replacing and:

- Add a first parameter corresponding to the receiver;
- Have the rest of the parameters match the method;
- Freely choose parameters to be encoded or decoded.

Therefore, during the COERCE phase, which introduces extension methods, the `implicit_*` has the signature:

```
def implicit_*(recv: @repr (Int, Int), n2:
  @repr (Int, Int)): @repr (Int, Int)
```

Since the programmer defining the description object is free to choose the encoding of the bypass arguments, the following (suboptimal) signature would be equally accepted:

```
def implicit_*(recv: (Int,Int), n2: (Int,Int)): (Int,Int)
```

It is interesting to notice that representation-specific method rewriting relies on two previous design choices: (1) shifting focus from scopes to individual values and (2) carrying the entire transformation semantics in the signature of each encoded value.

4.5 Interaction with Other Language Features

This section presents the interaction between the ADR transformation and object-oriented inheritance, generics and implicit conversions, explaining the additional steps that should be taken to ensure correct program transformation and the limitations of the approach.

Dynamic Dispatch and Overriding are an integral part of the object-oriented programming model, allowing objects to encapsulate code. The main approach to evolving this encapsulated code is extending the class and overriding its methods. However, changing the data representation can lead to situations where source-level overriding methods are no longer overriding in the low-level bytecode:

```
class X {
  def identity(i: (Int, Int)): (Int, Int) = i
}
adrt(IntPairAsLong) {
  class Y extends X(t: (Int, Int)) {
    override def identity(i: (Int, Int)) = t
  }
}
```

After the ADR transformation, the `identity` method in class `Y` no longer overrides method `identity` in class `X`, since its signature expects a long integer instead of a pair of integers. To address this problem, we extend the Late Data Layout mechanism, to introduce a new BRIDGE phase, which runs just before COERCE and inserts bridge methods to enable correct overriding. After the INJECT phase, the code corresponding to class `Y` is:

```
class Y extends X(t: @repr(...) (Int, Int)) {
  override def identity(i: @repr(...) (Int, Int)):
    @repr(...) (Int, Int) = t
}
```

The BRIDGE phase inserts the methods necessary to allow correct overriding (return types are omitted):

```
class Y extends X(t: @repr(...) (Int, Int)) {
  def identity(i: @repr(...) (Int, Int)) = t
  @bridge // overrides method identity from class X:
  override def identity(i: (Int, Int)) = identity(i)
}
```

The COERCE and COMMIT phases then transform class `Y` as before, resulting in a class with two methods, one containing the optimized code and another that overrides the method from class `X`, marked as `@bridge`:

```
class Y extends X(t: Long) {
  def identity(i: Long): Long = t
  @bridge override def identity(i: (Int, Int)) =
    IntPairAsLong.toHigh(identity(toRepr(i)))
}
```

If we now try to extend class `Y` in another `adrt` scope with the same transformation description object, overriding will take place correctly: the new class will define both the transformed method and the bridge, overriding both methods above. However, a more interesting case occurs when extending class `Y` from a scope with a different description:

```
adrt(IntPairAsDouble) { // != IntPairAsLong
  class Z(t: (Int, Int)) extends Y(t) {
    override def identity(i: (Int, Int)): (Int, Int) = i
  }
}
```

The ensuing BRIDGE phase generates 2 bridge methods:

```
class Z(t: Double) extends Y(...) {
  def identity(i: Double): Double = i
  @bridge override def identity(i: (Int, Int)) = ...
  @bridge override def identity(i: Long): Long = ...
}
```

Although the resulting object layout is consistent, the `@bridge` methods have to transform between the representations, which makes them less efficient. This is even more problematic when up-casting class `Z` to `Y` and invoking `identity`, as the bridge method goes through the high-level type to convert the long integer to a double. In such cases the BRIDGE phase issues warnings to notify the programmer of a possible slowdown caused by the coercions.

Generics. Another question that arises when performing ad hoc programmer-driven transformations is how to transform the data representation in generic containers. Should the ADR transformation be allowed to change the data representation stored in a `List`? We can use an example:

```
def use1(list: List[(Int, Int)]): Unit = ...
adrt(IntPairAsLong) {
  def use2(list: List[(Int, Int)]): Unit = use1(list)
}
```

In the specific case of the Scala immutable list, it would be possible to convert the `list` parameter of `use2` from type `List[Long]` to `List[(Int, Int)]` before calling `use1`. This can be done by mapping over the list and transforming the representation of each element. However, this domain-specific knowledge of how to transform the collection only applies to the immutable list in the standard library, and not to other generic classes that may occur in practice.

Furthermore, there is an entire class of containers for which this approach is incorrect: mutable containers. An

invariant of mutable containers is that any elements changed will be visible to all the code that holds a reference to the container. But duplicating the container itself and its elements (stored with a different representation) breaks this invariant: changes to one copy of the mutable container are not visible to its other copies. This is similar to the mutability restriction in §4.1.

The approach we follow in the ADR transformation is to preserve the high-level type inside generics. Thus, our example after the COMMIT phase will be:

```
1 def use1(list: List[(Int, Int)]): Unit = ...
2 def use2(list: List[(Int, Int)]): Unit = use1(list)
```

However, this does not prevent a programmer from defining another transformation description object that targets `List[(Int, Int)]` and replaces it by `List[Long]`:

```
1 def use1(list: List[(Int, Int)]): Unit = ...
2 adrt(ListOfIntPairAsListOfLong) {
3   def use2(list: List[(Int, Int)]): Unit = use1(list)
4 }
```

After the COMMIT phase, the transformation produces:

```
1 def use1(list: List[(Int, Int)]): Unit = ...
2 def use2(list: List[Long]): Unit = use1(toHigh(list))
```

Therefore `adrt` scopes are capable of targeting:

- generic types, such as `List[T]` for any `T`;
 - instantiated generic types, such as `List[(Int, Int)]`;
 - monomorphic types, such as `(Int, Int)`, outside generics
- Using these three cases and scope composition, programmers can conveniently target any type in their program.

Implicit conversions interact in two ways with `adrt` scopes: *Extending the object functionality* through implicit conversions, extension methods, or type classes must be taken into account by the method call rewriting in the COERCE phase. The handling of all three means of adding object functionality is similar, since, in all three cases, the call to the new method needs to be intercepted and redirected. Depending on the exact means, the mangled name for the bypass method will be different, but the mechanism and signatures remain the same (§4.4).

Offering an alternative to the LDL-based backend. Despite the apparent similarity, implicit conversions do not offer strong enough guarantees to replace the LDL backend. For example, assuming the presence of implicit methods to coerce integer pairs to longs and back, we can try to transform:

```
1 val n: (Int, Int) = (1, 0)
2 val a: Any = n
3 println(a)
```

To trigger the transformation, we replace the type of `n`:

```
1 val n: Long = implicitIntPairToLong((1, 0))
2 val a: Any = n // no coercion
3 println(a)
```

This resulting code changes semantics because no coercion is applied to `a`, since `Long` is a subtype of `Any`. In turn, this makes the output 4294967296 instead of `(1, 0)`. As we saw in §3, the missing coercion is correctly inserted when

annotations track the value representation, since annotations are orthogonal to the host language type system.

With this, we presented the main insights in the Ad hoc Data Representation Transformation approach and how they interact with other language features to guarantee transformation correctness. The next section describes the architecture and implementation of our Scala compiler plugin.

5. Implementation

This section describes the technical aspects of our implementation that can aid a compiler developer in porting the approach to another language. We implemented the ADR transformation as a Scala compiler plugin [2], by extending the open-source multi-stage programming transformation provided with the LDL [42] artifact, available at [4].

The `adrt` **scope** acts as the trigger for the ADR transformation. We treat it as a special keyword that we transform immediately after parsing, in the POSTPARSER phase. To show this, we follow a program through the compilation stages:

```
1 def foo: (Int, Int) = {
2   adrt(IntPairToLong) {
3     val n: (Int, Int) = (2, 4)
4   }
5   n
6 }
```

Immediately after the source is parsed, the POSTPARSER phase transforms the `adrt` scopes in three steps:

- it attaches a unique id to each `adrt` scope;
- it records and clears the block enclosed by the `adrt` scope
- it inlines the recorded code immediately after the now-empty `adrt` scope and, in the process, it marks the value and method definitions by the `adrt` scope's unique id (or by multiple ids, if `adrt` scopes are nested).

Following the POSTPARSER phase, the code is:

```
1 def foo: (Int, Int) = {
2   /* id: 100 */ adrt(IntPairToLong) {}
3   /* id: 100 */ val n: (Int, Int) = (2, 4)
4   n
5 }
```

This code is ready for type-checking: the definition of `n` is located in the same block as its use, making the scope correct. During the type-checking process, the `IntPairToLong` object is resolved to a symbol, missing type annotations are inferred and implicit conversions are introduced explicitly in the tree. After type-checking and pattern matching expansion, the INJECT phase traverses the tree and:

- for every `adrt` scope it records the id and description object, before removing it from the abstract syntax tree;
- for value and method definitions, if the type matches one or more transformations, it adds the `@repr` annotation.

Following the INJECT phase, the code for our example is:

```
1 def foo: (Int, Int) = {
2   val n: @repr(IntPairToLong) (Int, Int) = (2, 4)
3   n
4 }
```

After the INJECT phase, the annotated signatures are persisted, allowing the scope composition to work across sepa-

rate compilation. Later, the BRIDGE, COERCE and COMMIT phases proceed as described in §3 and §4.

The transformation description objects extend the marker trait `TransformationDescription`. Although the marker trait is empty, the description object needs to define at least the `toHigh` and `toRepr` coercions, which may be generic (§4.1). The programmer is then free to add bypass methods, in order to avoid decoding the representation type for the purpose of dynamically dispatching method calls. To aid the programmer in adding bypass methods, the COERCE phase warns whenever it does not find a suitable bypass method, indicating both the expected name and the expected method signature. Here we encountered a bootstrapping problem: although bypass methods handle the representation type, during the COERCE phase, their signatures are expected to accept parameters of the annotated high-level type, in order to allow redirecting method calls. To work around this problem, we added a the `@high` annotation, which acts as an anti-`@repr` and marks the representation types:

```
1 object IntPairToLong extends TransformationDescription{
2   ...
3   // source-level signature (type-checking the body):
4   def bypass_toString(repr: @high Long): String = ...
5   // signature during coerce (allows rewriting calls):
6   // def bypass_toString(repr: @repr(...)) (Int, Int)
7   // signature after commit (bytecode signature):
8   // def bypass_toString(repr: Long)
9 }
```

This mechanism allows programmers to both define and use the transformation description objects in the same compilation run. Considering the difficult nature of bootstrapping transformations, we are content with the current solution. Another advantage we get for free, thanks to referencing the transformation description object in the type annotation, is an explicit dependency between all transformed values and their description objects. This allows the Scala incremental compiler to automatically recompile all scopes when the description object in their `adrt` marker has changed.

This concludes the section, which explained how we solved the two main technical problems in the ADR Transformation and how this impacted the compilation pipeline.

6. Benchmarks

This section evaluates the experimental benefits of ADR transformations in targeted micro-benchmarks and in the setting of a library and its clients. We ran the benchmarks on an Intel i7-4702HQ quad-core processor machine with the frequency fixed at 2.2GHz. The RAM available to the benchmarks was 2GB and the running times were measured using the scalometer benchmarking platform [31], to avoid jitter from the garbage collector and just-in-time compiler.

6.1 ADRT Micro-Benchmarks

We chose representative micro-benchmarks in order to cover a wide range of transformations using the `adrt` scope:

- the greatest common divisor algorithm, presented in §2;
- least squares benchmark + deforestation [44];
- averaging sensor readings + array of struct;

- computing the first 10000 Hamming numbers.

The benchmarks were optimized using our implementation of the `adrt` scope at [2] and are described in detail on the website [3]. We will proceed to explain the transformation in each benchmark, but, due to space constraints, the full descriptions are only available on the website.

The Gaussian Greatest Common Divisor is the running example described in §2 and used throughout the paper. It is a numeric and CPU-bound benchmark, where the main slowdown is caused by heap allocations.

The `adrt` transformation surrounds the tail-recursive `gcd` method and optimizes the pair of integer representation of Gaussian integers by encoding them in long integers. From this point of view, we can think of the transformation as retrofitting the value class status to the pair of integers. The benchmark results in Table 1 show a 13x speed improvement and we checked that no object allocations occur in the transformed version of the `gcd` method. The transformation description object is 30 lines of code (LOC) long.

The Least Squares Method takes a list of points in two dimensions and computes the slope and offset of a straight line that best approximates the input data. The benchmark performs multiple traversals over the input data and thus benefits from deforestation [44], which avoids the creation of intermediate collections after each `map` operation:

```
1 adrt(ListAsLazyList){
2   def leastSquares(data: List[(Double, Double)]) = {
3     val size = data.length
4     val sumx = data.map(_._1).sum
5     val sumy = data.map(_._2).sum
6     val sumxy = data.map(p => p._1 * p._2).sum
7     val sumxx = data.map(p => p._1 * p._1).sum
8     ...
9   }
10 }
```

The `adrt` scope performs a generic transformation from `List[T]` to `LazyList[T]`:

```
1 object ListAsLazyList extends
2   TransformationDescription {
3   def toRepr[T](list: List[T]): LazyList[T] = ...
4   def toHigh[T](list: LazyList[T]): List[T] = ...
5   // bypass methods
6 }
```

The `LazyList` collection achieves deforestation by recording the mapped functions and executing them lazily, either when `force` is invoked on the collection or when a `fold` operation is executed. Since the `sum` operation is implemented as a `foldLeft`, the `LazyList` applies the function and sums the result without creating an intermediate collection. This transformation alone produced a 3x speedup for an input of 5 million points and made the `adrt` transformation perform better than the dedicated collection optimization tool `scalablitz` [7, 12], which produced a modest 1.7x speedup.

Yet, using the `LazyList` collection we can also benefit from specialization [19]. Using a specialized version of the `LazyList` collection we obtain a 5x speed improvement thanks to a combination of deforestation and specialization. Therefore, we used the `adrt` transformation both to transform the collection semantics and to retrofit specialization

in a localized scope. The transformation description object is 30 LOC and the `LazyList` is 70 LOC.

Although the Least Squares Method benchmark is a specialized micro-benchmark, the techniques it uses are of wide applicability. For instance, deforestation can be applied in a variety of realistic optimization scenarios.

The Sensor Readings benchmark was inspired by the Sparkle visualization tool [9], which is able to quickly display, zoom, transform and filter sensor readings. To obtain nearly real-time results, Sparkle combines several optimizations such as streaming and array-of-struct to struct-of-array conversions, all currently implemented by hand. In our benchmark, we implemented a mock-up of the processing core of Sparkle and automated the array-of-struct transform:

```
1 type SensorReadings = Array[(Long, Long, Double)]
2 class StructOfArray(arrayOfTimestamps: Array[Long],
3   arrayOfEvents: Array[Long],
4   arrayOfReadings: Array[Double])
5
6 object AoSToSoA extends TransformationDescription {
7   def toRepr(aos: SensorReadings): StructOfArray = ...
8   def toHigh(soa: StructOfArray): SensorReadings = ...
9   ...
10 }
```

Using the `adrt` scope produced a speedup of 2.2x. For comparison, we implemented the transformation by hand, manually refactoring the code, which obtained a nearly-identical speedup of 2.25x. The transformation description object is 60 LOC and the optimization applied (array-of-struct to struct-of-array) is general-purpose.

The Hamming Numbers Benchmark computes numbers that only have 2, 3 and 5 as their prime factors, in order. Unlike the other benchmarks, this is an example we randomly picked from Rosetta Code [6] and attempted to speed up:

```
1 adrt(BigIntToLong) {
2   adrt(QueueOfBigIntAsFunQueue) {
3     class Hamming extends Iterator[BigInt] {
4       import scala.collection.mutable.Queue
5       val q2 = new Queue[BigInt]
6       val q3 = new Queue[BigInt]
7       val q5 = new Queue[BigInt]
8       def enqueue(n: BigInt) = {
9         q2.enqueue(n * 2)
10        q3.enqueue(n * 3)
11        q5.enqueue(n * 5)
12      }
13       def next = {
14         val n = q2.head min q3.head min q5.head
15         if (q2.head == n) q2.dequeue
16         if (q3.head == n) q3.dequeue
17         if (q5.head == n) q5.dequeue
18         enqueue(n); n
19       }
20       def hasNext = true
21       q2.enqueue 1
22       q3.enqueue 1
23       q5.enqueue 1
24     }
25   }
26 }
```

An observation is that, for the first 10000 Hamming numbers, there is no need to use `BigInt`, since the numbers fit into a `Long` integer. Therefore, we used two nested `adrt` scopes to replace `BigInt` by `Long` and `Queue[BigInt]` by a fixed-size circular buffer built on an array. The result was

Benchmark	Original	ADRT	Speedup
Gaussian GCD	3.05 μ s	0.23 μ s	13x
Least Sq. Blitz (5M)	8026 ms	4763 ms	1.7x
Least Sq. adrt 1 (5M)	8026 ms	2393 ms	3x
Least Sq. adrt 2 (5M)	8026 ms	1643 ms	5x
Sensor Readings (5M)	50.8 ms	23.1 ms	2x
Hamming 10000th	4.35 ms	0.55 ms	8x

Table 1. Benchmark running time for each use case.

an 8x speedup. The main point in the transformation is its optimistic nature, which makes the assumption that, for the Hamming numbers we plan to extract, the long integer and a fixed-size circular buffer are good enough. This is similar to what a dynamic language virtual machine would do: it would make assumptions based on the code and would automatically de-specialize the code if the assumption is invalidated. In our case, when the assumption is invalidated, the code will throw an exception. For this example, the transformation has 100 LOC.

6.2 ADRT in Realistic Libraries

The `adrt` scoped transformation is a conceptual generalization of a mechanism motivated by library transformation scenarios. In particular, the resulting data representation transformation is used in conjunction with the miniboxing transformation [5, 41], in order to replace standard library *functions* and *tuples* by custom, optimized versions adequate for miniboxed code. The scope of this data representation transformation is miniboxing-transformed code.

The miniboxing transformation [41] proposes an alternative to erasure, allowing generic methods and classes to work efficiently with unboxed primitive types. Unlike the current specialization transformation in the Scala compiler [19], which duplicates and adapts the generic code once for every primitive type, the miniboxing transformation only duplicates the code once and *encodes all primitive types in long integers*. This allows miniboxing to scale much better than specialization [21] in terms of bytecode size while providing comparable performance. Yet, one of the main drawbacks of using the miniboxing plugin is that all Scala library classes are either generic or specialized with the built-in Scala specialization scheme, which is not compatible with miniboxing. Therefore, interacting with functions and tuples from miniboxed code incurs significant overhead.

Consider, for example, functions. (Tuples raise similar issues.) Scala offers functions as first-class citizens. However, since functions are not first-class citizens in the Java Virtual Machine bytecode, the Scala compiler desugars them to anonymous classes extending a functional interface. The following example shows the desugaring of function

`(x: Int) => x + 1:`

```
1 class $anon extends Function1[Int, Int] {
2   def apply(x: Int): Int = x + 1
3 }
4 new $anon()
```

This function desugaring does not expose a version of the `apply` method that encodes the primitive type as a long in-

Benchmark	Generic	Miniboxed	Miniboxed +functions
Sum	100.6 ms	355.9 ms	12.0 ms
SumOfSquares	188.3 ms	450.9 ms	13.0 ms
SumOfSqEven	130.8 ms	300.4 ms	52.2 ms
Cart	220.6 ms	560.2 ms	55.3 ms

Table 2. Scala Streams pipelines for 10M elements.

teger, as the miniboxing transformation expects. Therefore, when programmers write miniboxed code that uses functions, they have two choices: either accept the slowdown caused by converting the representation or define their own miniboxed `Function1` class, and perform the function desugaring by hand. Neither of these is a good solution.

Our data representation transformation converts the references to `Function1` in miniboxed code to the optimized `MiniboxedFunction1`, which allows calls to use the miniboxed representation, thus being more efficient. The problem is that the miniboxed code needs to interoperate with library-defined code, or with other libraries that were not transformed. Thus the miniboxed code acts as a scope for the *function and tuple representation transformation*, i.e., the ADR transformation of `Function` and `Tuple`. This transformation has a significant impact in library benchmarks.

The Scala-Streams library [13] imitates the design of the Java 8 stream library, to achieve high performance (relative to standard Scala libraries) for functional operations on data streams. The library is available as an open-source implementation [1]. In its continuation-based design, each stream combinator provides a function that is stacked to form a transformation pipeline. As the consumer reads from the final stream, the transformation pipeline is executed, processing an element from the source into an output element. However, the pipeline architecture is complex, since combinators such as `filter` may drop elements, stalling the pipeline.

Table 2 shows the result of applying our data representation transformation to the Scala-Streams published benchmarks. (The benchmarks are described in detail in prior literature [13].) As can be seen, the miniboxing transformation is an enabler of our optimization but produces *worse* results by itself (due to extra conversions).

Compared to the original library, the application of miniboxing and data representation optimization for functions achieves a very high speedup—up to 14.5x for the `SumOfSquares` benchmark. In fact, the speedup relative to the miniboxed code without the function representation optimization is nearly 35x!

The Framian Vector implementation is an exploration into deeply specializing the immutable `Vector` bulk storage without using reified types [10, 11]. This is a benchmark created by a commercial entity using the Scala programming language. Table 3 shows a 4.4x speed improvement when the function representation is optimized and shows the ADR-transformed function code lies within a 10% margin compared to the fully specialized and manually optimized code.

Benchmark	Running time
Manual C-like code	0.650 μ s
Miniboxing with functions	0.705 μ s
Miniboxing without functions	3.080 μ s
Generic	13.409 μ s

Table 3. Mapping a 1K vector.

7. Related Work

Changing data representations is a well-established and time-honored programming need. Techniques for removing abstraction barriers appear in the literature since the invention of high-level programming languages and often target low-level data representations. However, our technique is distinguished by its automatic determination of when data representations should be transformed, while giving the programmer control of how to perform this transformation and on which scope it is applicable.

As discussed earlier, the standard optimizations that are closest to our approach are value classes [8] and class specialization [19, 41]. These are optimizations with great practical value, and most modern languages have felt a need for them. For instance, specialization optimizations have recently been proposed for adoption in Java, with full VM support [22]. Rose has an analogous proposal for value classes [33, 34] in Java. Unlike our approach, all the above are whole-program data representation transformations and receive limited programmer input (e.g., a class annotation).

Virtual machine optimizations often also manage to produce efficient low-level representations through tracing [20] or inlining and escape analysis [18, 36]. Furthermore, modern VMs, such as V8, Truffle [46] and PyPy [14] attempt specialization and inference of optimized layouts. However, the ability to perform complex inferences dynamically is limited, and there is no way to draw domain-specific knowledge from the programmer. Generally VM optimizations are often successful at approaching the efficiency of a static language in a dynamic setting, but not successful in reliably exceeding it.

In terms of transformations, we already presented the Late Data Layout [42] mechanism in the Scala setting. Similar approaches, with different specifics in the extent of type system and customization support, have been applied to Haskell [24]. Foundational work exists for ML, with Leroy [27] presenting a transformation for unboxing objects, with the help of the type system. Later work extends [39] and generalizes [35] such transformations. In terms of runtime-dispatched generics, we refer to the work on Napier88 [29] and the TIL compiler [38] [23].

In the specific setting of data structure specialization, the CoCo approach [47] adaptively replaces uses of Java collections with optimized representations. CoCo has a similar high-level goal as our techniques, yet focuses explicitly on collections only. Approaches that only target a finite number of classes (data structure implementations) can be realized entirely in a library. An adaptive storage strategy for Python collections [15], for instance, switches representations once

collections become polymorphic or once they acquire many elements. The Scala Blitz optimizer uses macros to improve collection performance [7, 12].

Among mechanisms for extending an interface, such as extension methods, implicit conversions [30] and type classes [45] we can also mention views, which allow data abstraction and extraction through pattern matching [43].

Multi-stage programming [37] is another technique that optimizes the data representation. Its Scala implementation, dubbed lightweight modular staging and can both optimize and even re-target parts of a program to GPUs [16, 32]. Yet, multi-stage programming scopes are not accessible from outside, making it impossible to call a transformed method or read a transformed value. Instead, the transformation scope is closed and nothing is assumed to be part of the interface. Hopefully, this will be improved by techniques such as the Yin-Yang staging front-end [25], based on Scala macros [17]. Another type-directed transformation in the Scala compiler is the pickling framework [28], also based on macros. Instead of transforming the data representation in-place, pickler combinators create serialization code that allows can efficiently convert an object to a wide range of formats.

8. Conclusion

In this paper, we presented an intuitive interface over a safe and composable programmer-driven data representation transformation, where the composition works not only across source files but also across separate compilation runs. The transformation takes care of all the tedium involved in using a different representation, by automatically introducing coercions and bridge methods where necessary, and optimizing the code via extension methods. Benchmarking the resulting transformation shows significant performance improvements, with speedups between 1.9x and 14.5x.

References

- [1] Scala Streams Repository. URL <https://github.com/biboudis/scala-streams>.
- [2] ADRT Scala Compiler Plugin Source Code, . URL <https://github.com/EnuTheZhid/adrt-plugin>.
- [3] ADRT Scala Compiler Plugin Documentation, . URL <https://github.com/EnuTheZhid/adrt-plugin/wiki>.
- [4] LDL-based Staging Scala Compiler Plugin. URL <https://github.com/miniboxing/staging-plugin>.
- [5] The Miniboxing plugin website. URL <http://scala-miniboxing.org>.
- [6] Rosetta Code Website. URL <http://rosettacode.org>.
- [7] ScalaBlitz Optimizer. URL <https://scala-blitz.github.io/>.
- [8] Scala SIP-15: Value Classes. URL <http://docs.scala-lang.org/sips/completed/value-classes.html>.
- [9] Sparkle Tool. URL <https://github.com/mighdoll/sparkle>.
- [10] Optimistic Respecialization Attempts 1-5, . URL <http://io.pellucid.com/blog/optimistic-respecialization>.
- [11] Optimistic Respecialization Attempt 6, . URL <http://io.pellucid.com/blog/optimistic-respecialization-attempt-6>.
- [12] P. Aleksandar, D. Petrashko, and M. Odersky. Efficient Lock-Free Work-stealing Iterators for Data-Parallel Collections. In *PDP '15*. IEEE, 2015.
- [13] A. Biboudis, N. Palladinos, and Y. Smaragdakis. Clash of the Lambdas. *ICOLPS*, 2014.
- [14] C. F. Bolz, A. Cuni, M. Fijalkowski, and A. Rigo. Tracing the meta-level: PyPy's tracing JIT compiler. In *ICOLPS*, Genova, Italy, 2009. ACM.
- [15] C. F. Bolz, L. Diekmann, and L. Tratt. Storage Strategies for Collections in Dynamically Typed Languages. In *OOPSLA, OOPSLA '13*. ACM, 2013.
- [16] 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 *PACT*. IEEE Computer Society, 2011.
- [17] E. Burmako. Scala Macros: Let Our Powers Combine!: On How Rich Syntax and Static Types Work with Metaprogramming. In *SCALA*. ACM, 2013.
- [18] A. Deutsch. On Determining Lifetime and Aliasing of Dynamically Allocated Data in Higher-order Functional Specifications. In *POPL*. ACM, 1990. .
- [19] I. Dragos. *Compiling Scala for Performance*. PhD thesis, École Polytechnique Fédérale de Lausanne, 2010.
- [20] A. Gal. Trace-based Just-in-time Type Specialization for Dynamic Languages. In *PLDI*. ACM, 2009.
- [21] A. Genêt, V. Ureche, and M. Odersky. Improving the Performance of Scala Collections with Miniboxing (EPFL-REPORT-200245). Technical report, EPFL, 2014. URL <http://scala-miniboxing.org/>.
- [22] B. Goetz. State of the Specialization, 2014. URL <http://cr.openjdk.java.net/~briangoetz/valhalla/specialization.html>.
- [23] R. Harper and G. Morrisett. Compiling Polymorphism Using Intensional Type Analysis. In *PoPL*. ACM, 1995.
- [24] S. L. P. Jones and J. Launchbury. Unboxed Values as First Class Citizens in a Non-Strict Functional Language. In *Functional Programming Languages and Computer Architecture*. Springer, 1991.
- [25] V. Jovanovic, V. Nikolaev, N. D. Pham, V. Ureche, S. Stucki, C. Koch, and M. Odersky. Yin-Yang: Transparent Deep Embedding of DSLs. Technical report, EPFL, 2013.
- [26] B. W. Lampson and H. E. Sturgis. Reflections on an Operating System Design. *Commun. ACM*, 1976. ISSN 0001-0782.
- [27] X. Leroy. Unboxed Objects and Polymorphic Typing. In *PoPL*. ACM, 1992.
- [28] H. Miller, P. Haller, E. Burmako, and M. Odersky. Instant Pickles: Generating Object-oriented Pickler Combinators for Fast and Extensible Serialization. In *OOPSLA*. ACM, 2013.
- [29] R. Morrison, A. Dearle, R. C. H. Connor, and A. L. Brown. An Ad Hoc Approach to the Implementation of Polymorphism. *ACM TOPLAS*, 1991.
- [30] B. C. Oliveira, T. Schrijvers, W. Choi, W. Lee, and K. Yi. The Implicit Calculus: A New Foundation for Generic Programming. In *PLDI*. ACM, 2012.
- [31] A. Prokopec. ScalaMeter. URL <http://axel22.github.com/scalameter/>.
- [32] T. Rompf and M. Odersky. Lightweight Modular Staging: A Pragmatic Approach to Runtime Code Generation and Compiled DSLs. In *GPCE*, 2010. .
- [33] J. Rose. Value Types and Struct Tearing, . URL https://blogs.oracle.com/jrose/entry/value_types_and_struct_tearing.
- [34] J. Rose. Value Types in the VM, . URL https://blogs.oracle.com/jrose/entry/value_types_in_the_vm.
- [35] Z. Shao. Flexible Representation Analysis. In *ICFP*. ACM, 1997.
- [36] L. Stadler, T. Würthinger, and H. Mössenböck. Partial Escape Analysis and Scalar Replacement for Java. In *CGO*. ACM, 2014.
- [37] W. Taha. A Gentle Introduction to Multi-stage Programming. In *Domain-Specific Program Generation*, pages 30–50. Springer, 2004.
- [38] D. Tarditi, G. Morrisett, P. Cheng, C. Stone, R. Harper, and P. Lee. TIL: A Type-Directed Optimizing Compiler for ML. In *PLDI*. ACM, 1996.
- [39] P. J. Thiemann. Unboxed Values and Polymorphic Typing Revisited. In *Functional Programming Languages and Computer Architecture*. ACM, 1995.
- [40] V. Ureche. Additional Material for "Unifying Data Representation Transformations (EPFL-REPORT-200246)". Technical report, EPFL, 2014.
- [41] V. Ureche, C. Talau, and M. Odersky. Miniboxing: Improving the Speed to Code Size Tradeoff in Parametric Polymorphism Translations. In *OOPSLA*, 2013.
- [42] V. Ureche, E. Burmako, and M. Odersky. Late Data Layout: Unifying Data Representation Transformations. In *OOPSLA '14*. ACM, 2014.
- [43] P. Wadler. Views: A Way for Pattern Matching to Cohabit with Data Abstraction. In *PoPL*. ACM, 1987.
- [44] P. Wadler. Deforestation: Transforming Programs to Eliminate Trees. In *ESOP '88*. North-Holland Publishing Co., 1988.
- [45] P. Wadler and S. Blott. How to Make Ad-hoc Polymorphism Less Ad hoc. In *PoPL*. ACM, 1989.
- [46] T. Würthinger, C. Wimmer, A. Wöß, L. Stadler, G. Duboscq, C. Humer, G. Richards, D. Simon, and M. Wolczko. One VM to Rule Them All. In *Onward!* ACM, 2013.
- [47] G. Xu. CoCo: Sound and Adaptive Replacement of Java Collections. In *ECOOP*. Springer-Verlag, 2013.