

Additional Material for “Unifying Data Representation Transformations”

Vlad Ureche
vlad.ureche@epfl.ch

May 23, 2014

Abstract

This report shows an end-to-end formalization of the data representation transformation mechanism in the “Unifying Data Representation Transformations” paper [16]. Since the mechanism described in the paper is targeted at the Scala programming language and the specification is written against System $F_{<}$, with local colored type inference [11, 14] formally reasoning about the calculus is a major undertaking.

Instead, in this report we start from the simply typed lambda calculus with subtyping $\lambda_{<}$, natural numbers and unit. We add rewriting and adapt the calculus to propagate expected type information in a mechanism inspired from local colored type inference [11]. Finally we show how the representation transformation mechanism (the convert phase) rewrites terms. We prove that, given a series of assumptions about the inject phase, type-checking a term against the updated rules produces a correct and operationally equivalent term, with a minimum number of runtime coercions introduced for the annotations given.

We finish the report by giving a series of examples which show how the code is transformed.

1 Introduction

This report will present the formal aspects of the “Unifying Data Representation Transformations” paper [16]. Before diving into details, we will briefly refresh the main ideas in the paper.

1.1 Motivation

Programs can be seen as transformations between input and output data. Inside this transformation, processing the data involves passing values from one function to another, writing the data to variables, reading it back and performing primitive operations on it, such as addition and multiplication. In this context, the question of representation naturally arises: how does the program represent the data it manipulates. The question is particularly important when using high-level languages that isolate the programmer from the low-level processor architecture and the

realities of the platform the program is executed on. Inefficiently representing data has been shown to decrease performance and inflate the heap requirements [8, 6, 7, 4, 17]. Ideally, data should be represented in the most efficient format possible for the task. For example, when operating with integers, values should be stored in the processor registers or on the stack. Yet, certain high-level language features, such as parametric polymorphism (also called generics [3]) restrict the data representations that can be used in a low level program.

Parametric polymorphism allows programmers to abstract over the type of the data they are working with. This enables the definition of generic data structures, which operate in the same manner regardless of the data they are storing. For example, a linked list will perform the exact same lookup procedure regardless of whether the elements stored are integers or strings. This exposes a uniform interface to programmers, promoting code reuse while maintaining type safety. In the following example, written in the Scala programming language [10], the `identity` function takes any data type `T` and returns the exact same data type `T`:

```
1 def identity[T](arg: T): T = arg
2 val x: Int = identity[Int](0)
```

Yet, on the low level, different types of data use very different representations: integers are stored in general processor registers, floating point numbers use dedicated registers and references are represented as pointers to the heap memory. There is clearly a tension between the high-level uniformity and the low-level architecture specialization for each primitive type. In the simplest implementation of generics, the erasure transformation [3], data is represented uniformly both at the high level and at the low level: regardless of the type, the low level representation uses pointers to heap objects. This mandates that integers, which are best stored directly in processors registers, need to be represented as heap objects and passed by reference whenever interacting with generics. Under erasure, the previous example is translated to:

```
1 def identity(arg: Object): Object = arg
2 val x: Int = unbox(identity(box(0)))
```

Under erasure, the `box` and `unbox` data representation coercions are introduced as part of the translation. Yet, introducing coercions between data representations in an optimal fashion is a difficult task, which has been studied extensively [6, 8, 15]. Still, to this day, the algorithms for managing different data representations are custom tailored for specific tasks and do not have the means to express general representation constraints and to find the optimal points to introduce coercions such that the constraints are satisfied. The “Unifying Data Representation Transformation” paper [16] proposes a general and flexible mechanism for expressing data representation constraints in the type system and uses local type inference to automatically and optimally introduce coercions between different representations.

The general mechanism for data representation transformations uses annotated types [1, 2] to mark the representation constraints. This allows

local type inference to propagate the representation information along with the types. For example, in Scala, `Int` is considered the boxed representation while `@unboxed Int` is the unboxed integer. Explicitly propagating the representation in types allows the type system to automatically and optimally infer where coercions need to be introduced. This allows expressing a wide range of transformations in terms of the data representation transformation mechanism: (1) autoboxing, which automatically boxes and unboxes primitive types when necessary, as we have seen with generics, (2) value classes, which combine the advantages of flat structures with object-orientation, (3) specialization, which optimizes the translation of generics and (4) multi-stage programming, which allows evaluating the result of a program in multiple cycles containing run, generate optimized code, compile and launch phases.

The data representation transformation mechanism uniformly expresses all these use cases in as a three-step transformation:

- In the injection phase, transformations mark the types to use alternative representations. This is done using annotated types [1, 2];
- In the convert phase, type inference [13, 14, 11] and the annotations are used to guide optimal introduction of coercions between representations;
- Finally, having marked where and how the coercions take place, the data representation transformation can use the final semantics of the alternative data representations.

In the `identity` example, the autoboxing transformation would use the inject phase to mark values to be represented in their direct representation:

```
1 def identity[T](arg: T): T = arg
2 val x: @unboxed Int = identity[Int](0: @unboxed)
```

In the type application `identity[Int]`, the integer type is not marked as `@unboxed` since erased generics do require using boxed values. The next phase, convert, automatically introduces coercions between unboxed and boxed values, using the type system as a guide:

```
1 def identity[T](arg: T): T = arg
2 val x: @unboxed Int = unbox(identity[Int](box(0: @unboxed)))
```

Finally, in the commit phase, the autoboxing transformation can apply erasure and transform the code to its final form where `int` represents an unboxed integer:

```
1 def identity(arg: Object): Object = arg
2 val x: int = unbox(identity(box(0)))
```

The same mechanics can be employed by other data representation transformations as well, therefore making the mechanism an interesting object of study.

1.2 Formalization

In order to allow studying the data representation mechanism, it should have a formal specification that allows researchers to reason about it. Yet

the original paper does not include a formalization for several reasons:

- The complexity of the Scala type system and of System F with bounded quantification, System F_<;
- The domain-specific nature of the inject and commit phase, which are significantly different from a transformation to the next;
- The constraints of the paper, namely the page limit.

In order to provide an intuitive formalization, we restrict ourselves to the simply typed lambda calculus with subtyping $\lambda_{<}$, natural numbers and unit. To this end, we start by introducing a very simple local type inference-like mechanism to the calculus, which is done without interfering with its properties. Then we add annotated types to the language, prove progress and preservation. We then change the typing rules to account for the fact that `0`, `succ` and `pred` manipulate unboxed natural numbers. This point corresponds to the inject phase, where annotated types and their un-annotated counterparts are in a subtype relation in both directions, corresponding to a compatible state. Then, to formalize the convert phase, we remove the subtyping rules involving annotated types and introduce rewritings that add coercions. We prove by structural induction that any term that is typable in the first calculus is also typable in the second calculus, with the additional introduction of conversions. We then prove the rewritten term is operationally equivalent to the original term and prove that it produces programs that, on any given path, perform the minimum number of runtime coercions possible.

Throughout the formalization, we use the following example:

```
let identity : Nat → Nat = λarg : Nat. arg in
  let x : Nat = (identity 0) in
    unit
```

The `identity` function in the example is not generic, since the simply typed lambda calculus with subtyping $\lambda_{<}$ does not support generics. Instead we simulate generics by not marking the type of `arg` as `@unboxed`, therefore requiring the boxing and unboxing of the argument exactly as seen in the earlier example.

The next section will present the calculus.

2 The Calculus

This section will explain the calculus necessary for the data representation transformation mechanism, and will explain all the extensions necessary to express the injection and conversion phases. The calculus does not formalize the commit phase, which depends significantly on the data representations that are converted.

2.1 Simply Typed Lambda Calculus

We start with the simply typed lambda calculus with subtyping $\lambda_{<}$ and augment it with natural numbers in order to have primitive types which can be transformed. This is a clear step back from System F_< and the Scala type system, but it makes it possible to easily reason about the

$t ::=$ x $\lambda x : T. t$ $t t$ 0 $\text{succ } t$ $\text{pred } t$ unit	terms : <i>variable</i> <i>abstraction</i> <i>application</i> <i>constant zero</i> <i>successor</i> <i>predecessor</i> <i>unit</i>
$v ::=$ $\lambda x : T. t$ nv unit	values : <i>abstraction value</i> <i>numeric value</i> <i>unit</i>
$nv ::=$ 0 $\text{succ } nv$	numeric values : <i>zero value</i> <i>successor value</i>
$T ::=$ $T \rightarrow T$ Top Nat Unit	types : <i>type of functions</i> <i>maximum type</i> <i>natural numbers</i> <i>unit type</i>

Figure 1: Syntax of the $\lambda_{<}$ calculus with Nat and Unit

Desugaring:

$$t_1; t_2 \Rightarrow (\lambda x : \text{Unit} . t_2) t_1$$

where $x \notin FV(t_2)$

$$\text{let } x : T = e \text{ in } t \Rightarrow ((\lambda x : T . t) e)$$

Figure 2: Desugaring of the $\lambda_{<}$ calculus with Nat and Unit

Evaluation rules:

$$\frac{t_1 \longrightarrow t'_1}{t_1 t_2 \longrightarrow t'_1 t_2} \quad (\text{E-APP1})$$

$$\frac{t_2 \longrightarrow t'_2}{v_1 t_2 \longrightarrow v_1 t'_2} \quad (\text{E-APP2})$$

$$(\lambda x : T_1 . t_1) v_2 \longrightarrow [x \rightarrow v_2] t_1 \quad (\text{E-APPABS})$$

$$\text{pred } 0 \longrightarrow 0 \quad (\text{E-PREDZERO})$$

$$\text{pred } (\text{succ } nv_1) \longrightarrow nv_1 \quad (\text{E-PREDSUCC})$$

$$\frac{t_1 \longrightarrow t'_1}{\text{succ } t_1 \longrightarrow \text{succ } t'_1} \quad (\text{E-SUCC})$$

$$\frac{t_1 \longrightarrow t'_1}{\text{pred } t_1 \longrightarrow \text{pred } t'_1} \quad (\text{E-PRED})$$

Figure 3: Evaluation rules of the $\lambda_{<}$ calculus with `Nat` and `Unit`

transformation. We also augment the calculus with `unit`, to allow simple expression of the program.

One of the main missing features of this translation is parametric polymorphism, which arguably is the main motivation behind autoboxing and specialization. Yet both value classes and staging are use cases that are not motivated by generics. In fact staging is completely orthogonal to the low level implementation of generics in the language. In the motivating example, with the `identity` function, we instantiate the generic type to `Nat`:

```

let identity : Nat → Nat = λarg : Nat. arg in
  let x : Nat = (identity 0) in
    unit
    
```

To simulate generics in this case, the `inject` phase will not mark the argument type and the return type of `identity` for unboxing, therefore forcing the boxing and respectively unboxing of the argument and return type of the function.

On the other hand, we do include subtyping and the maximum type `Top`, which allows transforming interesting examples, such as:

```

1 val x : Object = identity(0)
    
```

Which would be represented as (omitting the definition of `identity`):

Typing rules:

$$\frac{x : T \in \Gamma}{\Gamma \vdash x : T} \quad (\text{T-VAR})$$

$$\frac{\Gamma, x : T_1 \vdash t_2 : T_2}{\Gamma \vdash (\lambda x : T_1. t_2) : T_1 \rightarrow T_2} \quad (\text{T-ABS})$$

$$\frac{\Gamma \vdash t_1 : T_1 \rightarrow T_2 \quad \Gamma \vdash t_2 : T_1}{\Gamma \vdash t_1 t_2 : T_2} \quad (\text{T-APP})$$

$$\Gamma \vdash 0 : \text{Nat} \quad (\text{T-ZERO})$$

$$\frac{\Gamma \vdash t_1 : \text{Nat}}{\Gamma \vdash \text{succ } t_1 : \text{Nat}} \quad (\text{T-SUCC})$$

$$\frac{\Gamma \vdash t_1 : \text{Nat}}{\Gamma \vdash \text{pred } t_1 : \text{Nat}} \quad (\text{T-PRED})$$

$$\frac{\Gamma \vdash t : S \quad S <: T}{\Gamma \vdash t : T} \quad (\text{T-SUB})$$

$$\Gamma \vdash \text{unit} : \text{Unit} \quad (\text{T-UNIT})$$

Figure 4: Typing rules of the $\lambda_{<}$ calculus with Nat and Unit

Subtyping rules:

$$S <: S \quad (\text{S-REFL})$$

$$\frac{S <: U \quad U <: T}{S <: T} \quad (\text{S-TRANS})$$

$$S <: \text{Top} \quad (\text{S-TOP})$$

$$\frac{T_1 <: S_1 \quad S_2 <: T_2}{S_1 \rightarrow S_2 <: T_1 \rightarrow T_2} \quad (\text{S-ARROW})$$

Figure 5: Subtyping rules of the $\lambda_{<}$ calculus with Nat and Unit

```
let x : Top = 0 in
  unit
```

The syntax of the calculus is given in Figure 1 and the desugaring in Figure 2. The operational semantics are given in Figure 3. Finally, the typing and subtyping rules are given in Figures 4 and 5.

For the $\lambda_{<}$ calculus with `Nat` and `Unit` we can rely on the inversion and substitution lemmas and on the progress and preservation theorems being proven in textbooks [13]:

2.2 Expected Type Propagation and Rewriting

Having shown the simply typed lambda calculus with subtyping, natural numbers and unit, we can now introduce the type propagation, which is the part we need from the local colored type inference and rewriting, which allows transforming the program during type checking. This includes adding two new elements:

- An **expected type propagation** mechanism inspired by the colored local type inference [11], which will later be used to propagate data representation constraints in the program. With this addition, type propagation becomes bi-directional: the syntax tree definitions have their type propagated inwards through the typing context Γ while expressions have their expected type propagated outwards by the expected type mechanism and
- A **term rewriting** extension, which is used to inject conversions at the right place, based on the divergence between the expected type and the actual type of a node. In the initial calculus in Figure 6, no term rewriting is performed: this is intentional, as we want to prove the properties of the calculus as we go. The further sections will add rewriting rules.

A natural question to ask is why not use Hindley-Milner [5, 9] full type reconstruction and use the more restricted expected type propagation technique. The key reason is that the expected type propagation mechanism allows rewriting terms as part of type checking, whereas the H-M type reconstruction only gathers constraints without updating the program. Later, if the representations do not match, this results in a set of irreconcilable constraints that prevent the type reconstruction algorithm from progressing, without clearly indicating where the constraint needs to be introduced. Another way to state this is that, using expected type propagation and rewriting, we allow the type system to rewrite program terms.

Another question that is worth answering is why not use implicit conversions which are available in the Scala type system. Indeed, implicit arguments have been formalized in [12]. Still, this state-of-the-art calculus does not formalize the introduction of implicit conversions based on the expected type. Instead, it focuses on implicit argument resolution, which does not require type propagation. This makes the implicit calculus in [12] impossible to use directly in both our formalization and in modelling implicit conversions. We hope that our work can serve as

the basis for extending the formalization to implicit conversions, since the convert phase could be seen as a simplified and restricted mechanism for implicit conversions. Still, from an implementation point of view, the convert phase cannot use the implicit resolution mechanism in the Scala compiler, since after name resolution import statements are removed from the syntax tree, therefore making it impossible to re-create the full scope necessary for resolving implicits.

With these two extensions, the typing judgement is the following:

$$T; \Gamma \vdash t \rightsquigarrow t' : T$$

We also introduce the wildcard type `*` in the type propagation scheme by using `*; \Gamma \vdash \dots`. The wildcard type is not a valid type in the program, but is used to mark the fact that no expected type has been propagated. In this case, `*; \Gamma \vdash t \rightsquigarrow t' : T` implies `T; \Gamma \vdash t \rightsquigarrow t' : T`. Contrarily, if a non-wildcard type is propagated, the rewritten member should conform to this expected type.

Figure 7 introduces the new typing rules. The syntax, desugaring and subtyping rules do not change. It is worth pointing out that the only way to satisfy an expected type in this version of the calculus is through the T-SUB rule, if the expression type-checks to a subtype of the expected type. This should explain why `*; \Gamma \vdash t \rightsquigarrow t' : T` implies `T; \Gamma \vdash t \rightsquigarrow t' : T`: by applying the T-SUB rule with the reflexive S-REFL subtyping rule on the left-hand side of the implication we immediately obtain the right-hand side.

Now that we have defined the simply typed lambda calculus with subtyping, natural numbers, unit, type propagation and rewriting, we can prove progress and preservation:

We will now add annotated types to the calculus.

2.3 Annotated Types

Annotated types allow marking the values to use an alternative representation. A very general solution is to allow annotating any type:

$T ::=$		$T \rightarrow T$	types :
		Top	<i>type of functions</i>
		Nat	<i>maximum type</i>
		Unit	<i>natural numbers</i>
		@unboxed T	<i>unit type</i>
			<i>annotated type</i>

This general annotation scheme is used in the Scala and Java type systems [1, 2]. Yet, the inject phase is more restrictive: it can only annotate those types that have an alternative representation. Let us assume that in our simple calculus, only natural number representation can be transformed. In this case, there is no semantic attached to type `@unboxed (Nat → Nat)`, since there is no alternative representation for functions. Instead, we can have `@unboxed Nat → @unboxed Nat`, which means a function from unboxed integers to unboxed integers. In order to build this restriction into the type system we can define types as:

Updated typing rules:

$$\frac{x : T \in \Gamma}{*, \Gamma \vdash x \rightsquigarrow x : T} \quad (\text{T-VAR})$$

$$\frac{*, \Gamma, x : T_1 \vdash t_2 \rightsquigarrow t'_2 : T_2}{*, \Gamma \vdash (\lambda x : T_1. t_2) \rightsquigarrow (\lambda x : T_1. t'_2) : T_1 \rightarrow T_2} \quad (\text{T-ABS})$$

$$\frac{*, \Gamma \vdash t_1 \rightsquigarrow t'_1 : T_1 \rightarrow T_2 \quad T_1; \Gamma \vdash t_2 \rightsquigarrow t'_2 : T_1}{*, \Gamma \vdash t_1 t_2 \rightsquigarrow t'_1 t'_2 : T_2} \quad (\text{T-APP})$$

$$*, \Gamma \vdash 0 \rightsquigarrow 0 : \text{Nat} \quad (\text{T-ZERO})$$

$$\frac{\text{Nat}; \Gamma \vdash t_1 \rightsquigarrow t'_1 : \text{Nat}}{*, \Gamma \vdash \text{succ } t_1 \rightsquigarrow \text{succ } t'_1 : \text{Nat}} \quad (\text{T-SUCC})$$

$$\frac{\text{Nat}; \Gamma \vdash t_1 \rightsquigarrow t'_1 : \text{Nat}}{*, \Gamma \vdash \text{pred } t_1 \rightsquigarrow \text{pred } t'_1 : \text{Nat}} \quad (\text{T-PRED})$$

$$\frac{*, \Gamma \vdash t \rightsquigarrow t' : S \quad S <: T}{T; \Gamma \vdash t \rightsquigarrow t' : T} \quad (\text{T-SUB})$$

$$*, \Gamma \vdash \text{unit} \rightsquigarrow \text{unit} : \text{Unit} \quad (\text{T-UNIT})$$

Figure 6: Typing rules of the $\lambda_{<}$ calculus with Nat, Unit, expected type propagation and rewriting

Updated subtyping rules:

$$\begin{array}{l}
 S <: S \qquad\qquad\qquad (\text{S-REFL}) \\
 \\
 \frac{S <: U \quad U <: T}{S <: T} \qquad\qquad\qquad (\text{S-TRANS}) \\
 \\
 \text{T1} <: \text{Top} \text{ only for types of T1, not annotated types} \quad (\text{S-TOP}) \\
 \\
 \frac{T_1 <: S_1 \quad S_2 <: T_2}{S_1 \rightarrow S_2 <: T_1 \rightarrow T_2} \qquad\qquad\qquad (\text{S-ARROW}) \\
 \\
 @\text{unboxed Nat} <: \text{Nat} \qquad\qquad\qquad (\text{S-NAT1}) \\
 \\
 \text{Nat} <: @\text{unboxed Nat} \qquad\qquad\qquad (\text{S-NAT2})
 \end{array}$$

Figure 7: Subtyping rules of the $\lambda_{<}$ calculus with Nat, Unit, expected type propagation, rewriting and annotations

T ::=		@unboxed Nat	all types : <i>unboxed natural numbers</i>
		T1	<i>unannotated types</i>
T1 ::=		$T \rightarrow T$	unannotated types : <i>type of functions</i>
		Top	<i>maximum type</i>
		Nat	<i>natural numbers</i>
		Unit	<i>unit type</i>

It is important to notice that, unlike the Scala type system, the simplified set of types above does not allow annotating a type twice. While for our use case this is the natural choice, both the type systems of Java and Scala allow multiple annotations on a type.

In this calculus, it is not yet necessary to separate Top and annotated types. Still, we can already do this, since it is required by the data representation mechanism anyway. To do so, we give the subtyping rules in Figure 7, where Top is only a supertype of any T1 but not of any T. The gap is filled by combining the S-NAT1 rule with S-TOP: $@\text{unboxed Nat} <: \text{Nat} <: \text{Top}$.

The next section will present the data representation mechanism.

3 The Data Representation Mechanism

Until now, we have built the calculus necessary to express the data representation transformation. Now we can finally express the transformations

Updated typing rules:

$$\begin{array}{l}
 *; \Gamma \vdash 0 \rightsquigarrow 0 : @unboxed \text{ Nat} \quad (\text{T-ZERO}) \\
 \\
 \frac{@unboxed \text{ Nat}; \Gamma \vdash t_1 \rightsquigarrow t'_1 : @unboxed \text{ Nat}}{*; \Gamma \vdash \text{succ } t_1 \rightsquigarrow \text{succ } t'_1 : @unboxed \text{ Nat}} \quad (\text{T-SUCC}) \\
 \\
 \frac{@unboxed \text{ Nat}; \Gamma \vdash t_1 \rightsquigarrow t'_1 : @unboxed \text{ Nat}}{*; \Gamma \vdash \text{pred } t_1 \rightsquigarrow \text{pred } t'_1 : @unboxed \text{ Nat}} \quad (\text{T-PRED})
 \end{array}$$

Figure 8: Updated typing rules for the injection phase

themselves.

3.1 The Inject Phase

The inject phase will inject the `@unboxed` annotation where necessary, following a custom logic. During the injection phase, the `succ`, `pred` and `0` nodes are also converted to `@unboxed Nat` instead of `Nat`, to signal the fact that their results are unboxed. The updated rules for the calculus are given in Figure 8.

With these changes, we can re-state the theorems:
Looking at our earlier example:

```

let identity : Nat → Nat = λarg: Nat. arg in
  let x : Nat = (identity 0) in
    let y : Top = 0 in
      unit
    
```

After the inject transformation it will look like:

```

let identity : Nat → Nat = λarg: Nat. arg in
  let x : @unboxed Nat = (identity 0) in
    let y : Top = 0 in
      unit
    
```

So far, the example has not changed much, a single annotation was added for `x`. The convert phase does the heavy lifting of adding explicit representation conversions.

3.2 The Convert Phase

The convert phase is the most complex phase in the transformation. It is the first phase to make full use of both the expected type propagation and the term rewriting capacity in the calculus we defined.

The most important element in the convert phase are the conversion nodes, which are introduced while rewriting the tree. Along with the conversion nodes, the values in the calculus need to be patched to account for boxed values, as shown if Figure 9. Now that explicit conversions

t	$::=$ x $\lambda x : T. t$ $t t$ 0 $\text{succ } t$ $\text{pred } t$ $\text{box } t$ $\text{unbox } t$ unit	terms : <i>variable</i> <i>abstraction</i> <i>application</i> <i>constant zero</i> <i>successor</i> <i>predecessor</i> <i>box number</i> <i>unbox number</i> <i>unit</i>
v	$::=$ $\lambda x : T. t$ nv $\text{box } nv$ unit	values : <i>abstraction value</i> <i>numeric value</i> <i>box numeric value</i> <i>unit</i>
nv	$::=$ 0 $\text{succ } nv$	numeric values : <i>zero value</i> <i>successor value</i>
T	$::=$ $@\text{unboxed Nat}$ $T1$	all types : <i>unboxed natural numbers</i> <i>unannotated types</i>
$T1$	$::=$ $T \rightarrow T$ Top Nat Unit	unannotated types : <i>type of functions</i> <i>maximum type</i> <i>natural numbers</i> <i>unit type</i>

Figure 9: Full syntax of the convert phase

Evaluation rules:

$$\frac{t_1 \longrightarrow t'_1}{t_1 t_2 \longrightarrow t'_1 t_2} \quad (\text{E-APP1})$$

$$\frac{t_2 \longrightarrow t'_2}{v_1 t_2 \longrightarrow v_1 t'_2} \quad (\text{E-APP2})$$

$$(\lambda x : T_1 . t_1) v_2 \longrightarrow [x \rightarrow v_2] t_1 \quad (\text{E-APPABS})$$

$$\text{pred } 0 \longrightarrow 0 \quad (\text{E-PREDZERO})$$

$$\text{pred } (\text{succ } nv_1) \longrightarrow nv_1 \quad (\text{E-PREDSUCC})$$

$$\frac{t_1 \longrightarrow t'_1}{\text{succ } t_1 \longrightarrow \text{succ } t'_1} \quad (\text{E-SUCC})$$

$$\frac{t_1 \longrightarrow t'_1}{\text{pred } t_1 \longrightarrow \text{pred } t'_1} \quad (\text{E-PRED})$$

$$\frac{t_1 \longrightarrow t'_1}{\text{box } t_1 \longrightarrow \text{box } t'_1} \quad (\text{E-BOX})$$

$$\text{box } (\text{unbox } t) \longrightarrow t \quad (\text{E-BOXUNBOX})$$

$$\frac{t_1 \longrightarrow t'_1}{\text{unbox } t_1 \longrightarrow \text{unbox } t'_1} \quad (\text{E-UNBOX})$$

$$\text{unbox } (\text{box } t) \longrightarrow t \quad (\text{E-UNBOXBOX})$$

Figure 10: Full evaluation rules of the convert phase

Updated subtyping rules for the convert phase:

$$\begin{array}{l}
 S <: S \qquad\qquad\qquad (\text{S-REFL}) \\
 \\
 \frac{S <: U \quad U <: T}{S <: T} \qquad\qquad\qquad (\text{S-TRANS}) \\
 \\
 T1 <: \text{Top only for types of T1, not annotated types} \quad (\text{S-TOP}) \\
 \\
 \frac{T_1 <: S_1 \quad S_2 <: T_2}{S_1 \rightarrow S_2 <: T_1 \rightarrow T_2} \qquad\qquad\qquad (\text{S-ARROW})
 \end{array}$$

Figure 11: Subtyping rules of the convert phase

have been introduced, they must also be considered into the operational semantics, as shown in Figure 10.

To trigger the introduction of conversions into the tree, we transform representation inconsistencies into mismatching types. The first step in doing so is to eliminate the subtyping relations S-NAT1 and S-NAT2 (introduced for the inject phase in Figure 7), which ensure compatibility between representations during the inject phase. The second step is to add two new typing rules T-ADAPTBOX and T-ADAPTUNBOX that explicitly introduce the `box` and `unbox` operations, in Figure 12. Finally, the T-BOX and T-UNBOX operations are added to allow re-type-checking the tree after the `box` and `unbox` operations were introduced.

Following the rewriting in the convert phase, our example will become:

```

let identity : Nat → Nat = λarg : Nat. arg in
  let x : @unboxed Nat = unbox (identity (box 0)) in
    let y : Top = box 0 in
      unit
    
```

In the initial example, we started from a state where `Nat` was the only numeric type. In the injection phase, we introduced the notion of `@unboxed Nat` and added `@unboxed` annotation to `x`. With this, the conversion phase added the explicit conversions necessary to have consistent data representations in the program. From a technical perspective, a later commit phase allows the compiler to give the final semantics of the `@unboxed` annotation, but since this phase is very different for each data transformation, we will not attempt to formalize it.

It must be noted that optimality is not about the minimum total number of conversions in the program. Instead, it refers to the minimum number of conversions taken on any execution path (trace) through the program, modulo the constraints introduced by the inject phase. Assuming we add booleans to our calculus with the `if` expression, we can take the following example:

Updated typing rules for the convert phase:

$$\begin{array}{c}
 \frac{x : T \in \Gamma}{*; \Gamma \vdash x \rightsquigarrow x : T} \quad (\text{T-VAR}) \\
 \\
 \frac{*; \Gamma, x : T_1 \vdash t_2 \rightsquigarrow t'_2 : T_2}{*; \Gamma \vdash (\lambda x : T_1. t_2) \rightsquigarrow (\lambda x : T_1. t'_2) : T_1 \rightarrow T_2} \quad (\text{T-ABS}) \\
 \\
 \frac{*; \Gamma \vdash t_1 \rightsquigarrow t'_1 : T_1 \rightarrow T_2 \quad T_1; \Gamma \vdash t_2 \rightsquigarrow t'_2 : T_1}{*; \Gamma \vdash t_1 t_2 \rightsquigarrow t'_1 t'_2 : T_2} \quad (\text{T-APP}) \\
 \\
 *; \Gamma \vdash 0 \rightsquigarrow 0 : @\text{unboxed Nat} \quad (\text{T-ZERO}) \\
 \\
 \frac{@\text{unboxed Nat}; \Gamma \vdash t_1 \rightsquigarrow t'_1 : @\text{unboxed Nat}}{*; \Gamma \vdash \text{succ } t_1 \rightsquigarrow \text{succ } t'_1 : @\text{unboxed Nat}} \quad (\text{T-SUCC}) \\
 \\
 \frac{@\text{unboxed Nat}; \Gamma \vdash t_1 \rightsquigarrow t'_1 : @\text{unboxed Nat}}{*; \Gamma \vdash \text{pred } t_1 \rightsquigarrow \text{pred } t'_1 : @\text{unboxed Nat}} \quad (\text{T-PRED}) \\
 \\
 \frac{*; \Gamma \vdash t \rightsquigarrow t' : S \quad S <: T}{T; \Gamma \vdash t \rightsquigarrow t' : T} \quad (\text{T-SUB}) \\
 \\
 \frac{*; \Gamma \vdash t \rightsquigarrow t' : @\text{unboxed Nat}}{\text{Nat}; \Gamma \vdash t \rightsquigarrow \text{box } t' : \text{Nat}} \quad (\text{T-ADAPTBOX}) \\
 \\
 \frac{*; \Gamma \vdash t \rightsquigarrow t' : \text{Nat}}{@\text{unboxed Nat}; \Gamma \vdash t \rightsquigarrow \text{unbox } t' : @\text{unboxed Nat}} \quad (\text{T-ADAPTUNBOX}) \\
 \\
 \frac{*; \Gamma \vdash t \rightsquigarrow t' : @\text{unboxed Nat}}{*; \Gamma \vdash \text{box } t \rightsquigarrow \text{box } t' : \text{Nat}} \quad (\text{T-BOX}) \\
 \\
 \frac{*; \Gamma \vdash t \rightsquigarrow t' : \text{Nat}}{*; \Gamma \vdash \text{unbox } t \rightsquigarrow \text{unbox } t' : @\text{unboxed Nat}} \quad (\text{T-UNBOX}) \\
 \\
 *; \Gamma \vdash \text{unit} \rightsquigarrow \text{unit} : \text{Unit} \quad (\text{T-UNIT})
 \end{array}$$

Figure 12: Typing rules of the convert phase


```
let pick : Bool → @unboxed Nat → @unboxed Nat → Top =
  λcond: Bool.
  λchoice1: @unboxed Nat.
  λchoice2: @unboxed Nat.
    if cond then
      box choice1
    else
      box choice2
  in
  unit
```

In this example, the transformation pushes the coercions as deep as possible in the tree, speculating that some of the branches will not need the coercion. Still, this is suboptimal from a global coercion count, since the translation could have been place around the if statement to have a single global coercion instead of tow. Yet, coercing each branch separately is more optimal in the following example:

```
let pick : Bool → Nat → @unboxed Nat → Top =
  λcond: Bool.
  λchoice1: Nat.
  λchoice2: @unboxed Nat.
    if cond then
      choice1 ; no coercion!
    else
      box choice2
  in
  unit
```

In the second example, the if $cond = \text{true}$ there are no coercions in the execution trace and if $cond = \text{false}$ there is one. One could even argue that, in case the if statement is duplicated, so is the coercion. This is true, but it is captured by the requirement that the annotations injected are satisfied – if the high-level annotations are suboptimal, as we have shown in the paper, so will the transformed program.

Intuitively, optimality comes from the fact that the convert phase will propagate expected types and will not filp the representation without a constraint. Since constraints are produced by interacting with code annotated in the inject phase, we can state that conversions are optimal modulo the constraints introduced by the inject phase.

The reminder of this document has hand-written proofs attached.

①

Proofs for the "Unifying
Data Representation Transformations"
Paper

Vlad Ureche

The following proofs describe both the inject (§3.1) and the convert phase (§3.2). Any difference is noted both in the theorem statement and in the theorem proof.

Inversion lemma

- ① If $*; \Gamma \vdash r \rightsquigarrow x : R$ then $r = x$ and $x : R \in \Gamma$
- ② If $*; \Gamma \vdash r \rightsquigarrow \lambda x : T_1. t_2' : R$ then $\exists T_2, t_2$ s.t. $R = T_1 \rightarrow T_2$
and $r = \lambda x : T_1. t_2$
and $*; \bar{x} : T_1, \Gamma \vdash t_2 \rightsquigarrow t_2' : T_2$
- ③ If $*; \Gamma \vdash r \rightsquigarrow t_1' t_2' : R$ then $\exists T_1, t_1, t_2$ s.t. $r = t_1 t_2$
and $*; \Gamma \vdash t_1 \rightsquigarrow t_1' : T_1 \rightarrow R$
Careful! \rightarrow and $\underline{T_1}, \Gamma \vdash t_2 \rightsquigarrow t_2' : T_1$
- ④ If $*; \Gamma \vdash r \rightsquigarrow 0 : R$ then $r = 0$ and $R = @ \text{unboxed Nat}$
- ⑤ If $*; \Gamma \vdash r \rightsquigarrow \text{succ } t_1' : R$ then $\exists t_1$ s.t. $r = \text{succ } t_1$
and $R = @ \text{unboxed Nat}$
and $@ \text{u Nat}; \Gamma \vdash t_1 \rightsquigarrow t_1' : @ \text{u Nat}$
- ⑥ If $*; \Gamma \vdash r \rightsquigarrow \text{pred } t_1' : R$ then $\exists t_1$ s.t. $r = \text{pred } t_1$
and $R = @ \text{unboxed Nat}$
and $@ \text{u Nat}; \Gamma \vdash t_1 \rightsquigarrow t_1' : @ \text{u Nat}$
- ⑦ If $*; \Gamma \vdash r \rightsquigarrow \text{unit} : R$ then $r = \text{unit}$
and $R = \text{Unit}$

⑧ If $T; \Gamma \vdash r \rightsquigarrow t' : R$ then $\exists s. s.t. s <: R. s.t. \textcircled{2}$
 $*; \Gamma \vdash r \rightsquigarrow t' : s.$

for the coerce phase:

⑨ If $*; \Gamma \vdash r \rightsquigarrow \text{box } t' : R$ then $\exists t. s.t. R = \text{Nat}$
 and $r = \text{box } t$
 and $*; \Gamma \vdash t \rightsquigarrow t' : @\text{unboxed Nat}$

⑩ If $*; \Gamma \vdash r \rightsquigarrow \text{unbox } t' : R$ then $\exists t. s.t. R = @\text{unboxed Nat}$
 and $r = \text{unbox } t$
 and $*; \Gamma \vdash t \rightsquigarrow t' : \text{Nat}$

⑪ If $T; \Gamma \vdash r \rightsquigarrow \text{box } t' : R$ then $\exists t. s.t. R = \text{Nat}$
 and $T = \text{Nat}$
 and $r = t$
 and $*; \Gamma \vdash t \rightsquigarrow t' : @\text{unboxed Nat}$

⑫ If $T; \Gamma \vdash r \rightsquigarrow \text{unbox } t' : R$ then $\exists t. s.t. R = @\text{unboxed Nat}$
 and $T = @\text{unboxed Nat}$
 and $r = t$
 and $*; \Gamma \vdash t \rightsquigarrow t' : \text{Nat}$

Theorem 1 [Uniqueness of Types And Rewritings]

Each term t has at most one type ^{and one rewriting}. That is, if t is typeable, then its type is unique. For injection only: Moreover, there is only one derivation of this typing built from the inference rules. For convert: We can have multiple typing derivations, but we will prove a weaker property, namely the idempotence property.

A mathematical specification is:
 $\exists z. \Gamma \vdash t \rightsquigarrow t' : T$ and $\exists z. \Gamma \vdash t \rightsquigarrow t'' : T$ OR
 $\exists z. \Gamma \vdash t \rightsquigarrow t' : T$ and $\exists z. \Gamma \vdash t \rightsquigarrow t'' : T$ then $T = T'$
and
 $t' = t''$

Proof Structural induction based on the relation that:

$$*; \Gamma \vdash t \rightsquigarrow t' : T$$

For inject: Lemma: The rewrites produce the same terms $\Rightarrow t = t'$

For convert: uniqueness of type and rewriting

- ① $*; \Gamma \vdash t \rightsquigarrow x : T \xRightarrow{\text{inv.}} t = x \text{ and } x : T \in \Gamma$ ③
 if $*; \Gamma \vdash x \rightsquigarrow t'' : T' \xRightarrow{T\text{-VAR}} t'' = x \text{ and } T' = T$
- ② $*; \Gamma \vdash t \rightsquigarrow \lambda x : T_1. t_2' : T \xRightarrow{\text{inv.}} \exists T_2, t_2 \text{ s.t. } T = T_1 \rightarrow T_2$
 and $t = \lambda x : T_1. t_2$
 and $*; x : T_1, \Gamma \vdash t_2 \rightsquigarrow t_2' : T_2$ (1)
 iH. $\Rightarrow T_2$ and t_2 are unique
 if $*; \Gamma \vdash \lambda x : T_1. t_2 \rightsquigarrow t'' : T' \xrightarrow{T\text{-ABS}} t'' = \lambda x : T_1. t_2'$
 $T' = T_1 \rightarrow T_2$
- ③ $*; \Gamma \vdash t \rightsquigarrow t_1' t_2' : T \xRightarrow{\text{inv.}} \exists T_1, t_1, t_2 \text{ s.t. } t = t_1 t_2$
 and $*; \Gamma \vdash t_1 \rightsquigarrow t_1' : T_1 \rightarrow T$ (2)
Careful! \rightarrow and $*; \Gamma \vdash t_2 \rightsquigarrow t_2' : T_1$ (3)
 by the induction hypothesis (IH) \Rightarrow
 t_1' and t_2' are unique. (4)
 $T_1 \rightarrow T$ and T_1 are unique (5)
 if $*; \Gamma \vdash t_1 t_2 \rightsquigarrow t'' : T' \xrightarrow[T=5]{T\text{-APP}} t'' = t_1' t_2'$ and $T' = T$
- ④ $*; \Gamma \vdash t \rightsquigarrow 0 : T \xRightarrow{\text{inv.}} t = 0 \text{ and } T = @ \text{ unboxed Nat}$
 if $*; \Gamma \vdash 0 \rightsquigarrow t'' : T' \xRightarrow{T\text{-ZERO}} t'' = 0 \text{ and } T' = T = @ \text{ unboxed Nat}$
- ⑤ $*; \Gamma \vdash t \rightsquigarrow \text{succ } t_1' : T \xRightarrow{\text{inv.}} \exists t_1 \text{ s.t. } t = \text{succ } t_1$
 and $T = @ \text{ unboxed Nat}$
 and $@ \text{ unboxed Nat}; \Gamma \vdash t_1 \rightsquigarrow t_1' : @ \text{ unboxed Nat}$
 by the IH $\Rightarrow t_1'$ is unique, can't type
 to any other type. (6)
 if $*; \Gamma \vdash \text{succ } t_1 \rightsquigarrow t'' : T' \xrightarrow[T=6]{T\text{-SUCC}} t'' = \text{succ } t_1'$
 $T' = T = @ \text{ unboxed Nat}$
- ⑥ pred \rightarrow identical. to succ
- ⑦ unit \rightarrow identical to 0

⑧ $\mathbb{Z}; \Gamma \vdash t \rightsquigarrow t' : T$ ④

(also covers cases of the inversion lemma)
 in this case, we actually have 3 sub-cases:
 (8.1) $\exists s, t. S <: \mathbb{Z} \text{ s.t.}$

$$\left. \begin{array}{l} *; \Gamma \vdash t \rightsquigarrow t' : S \\ \text{if } *; \Gamma \vdash t \rightsquigarrow t'' : S' \end{array} \right\} \stackrel{\text{IH}}{\Rightarrow} t' = t'' \text{ and } S = S' \quad (7)$$

and $T = \mathbb{Z}$.
 if $\mathbb{Z}; \Gamma \vdash t \rightsquigarrow t'' : T'$

only for convert, since in inject the only rule for expected types is T-Sub.
 Q: can another case apply?
 A: No, if we try to apply T-ADAPT Box, it means $\mathbb{Z} = \text{Nat}$ and $S = @\text{unboxed Nat}$, but $S <: \mathbb{Z}$ is not satisfied, since in convert we eliminated S-NAT1 and S-NAT2. Same reasoning for T-ADAPT UnBox.

$\stackrel{T\text{-Sub}}{\Rightarrow} t'' = t' \text{ and } T = T'$

(8.2) $t' = \text{box } t_0$, by the inversion lemma $\Rightarrow T = \mathbb{Z} = \text{Nat}$
 only for convert and $*; \Gamma \vdash t \rightsquigarrow t_0' : @\text{unboxed Nat}$ (8)

if $\text{Nat}; \Gamma \vdash t \rightsquigarrow t'' : T'' \stackrel{T\text{-ADAPT Box}}{\Rightarrow} t'' = \text{box } t_0' = t'$

$T'' = \text{Nat} = T$

* Q: is T-Sub applicable?

A: No, since $@\text{unboxed Nat} \not<: \text{Nat}$

(8.3) similar to (8.2)

only for convert

⑨ if $*; \Gamma \vdash t \rightsquigarrow \text{box } t_0' : T \stackrel{\text{inv}}{\Rightarrow} \exists t_0 \text{ s.t. } t = \text{box } t_0 \text{ and } T = \text{Nat}$

and $*; \Gamma \vdash t_0 \rightsquigarrow t_0' : @\text{unboxed Nat}$ } $\stackrel{\text{IH}}{\Rightarrow} t_0'' = t_0'$ and (9)
 $*; \Gamma \vdash t_0 \rightsquigarrow t_0'' : T'' \quad T'' = @\text{unboxed Nat}$

$*; \Gamma \vdash \text{box } t_0 \rightsquigarrow t'' : T' \stackrel{T\text{-Box}}{\Rightarrow} t'' = \text{box } t_0' = t'$
 (9) and $T' = T = \text{Nat}$

⑩ if $*; \Gamma \vdash t \rightsquigarrow \text{unbox } t_0' : T \Rightarrow \dots$ exactly the same as ⑨

\Rightarrow we have covered all cases ■

4 Conclusion

This report has shown the formalization of the “Unifying Data Representation Transformations” paper and aims at providing a base for other formalizations as well, especially in the area of implicit conversions (views).

References

- [1] JSR 308: Annotations on Java Types.
- [2] SIP-5 - Internals of Scala Annotations.
- [3] BRACHA, G., ODERSKY, M., STOUTAMIRE, D., AND WADLER, P. Making the future safe for the past: Adding Genericity to the Java Programming Language. In *OOPSLA* (1998), ACM.
- [4] DRAGOS, I., AND ODERSKY, M. Compiling Generics Through User-Directed Type Specialization. In *ICOOOLPS* (Genova, Italy, 2009).
- [5] HINDLEY, R. The Principal Type-scheme of an Object in Combinatory Logic. *Transactions of the American Mathematical Society* (1969).
- [6] JONES, S. L. P., AND LAUNCHBURY, J. Unboxed Values as First Class Citizens in a Non-Strict Functional Language. In *Functional Programming Languages and Computer Architecture* (1991), Springer.
- [7] KENNEDY, A., AND SYME, D. Design and Implementation of Generics for the .NET Common Language Runtime. In *PLDI* (2001).
- [8] LEROY, X. Unboxed Objects and Polymorphic Typing. In *PoPL* (1992), ACM.
- [9] MILNER, R. A Theory of Type Polymorphism in Programming. *Journal of Computer and System Sciences* (1978).
- [10] ODERSKY, M., SPOON, L., AND VENNERS, B. *Programming in Scala: A Comprehensive Step-by-step Guide*. Artima Inc, 2008.
- [11] ODERSKY, M., ZENGER, M., AND ZENGER, C. Colored Local Type Inference. In *PoPL* (2001), ACM.
- [12] OLIVEIRA, B. C., SCHRIJVERS, T., CHOI, W., LEE, W., AND YI, K. The Implicit Calculus: A New Foundation for Generic Programming. In *PLDI* (2012), ACM.
- [13] PIERCE, B. C. *Types and Programming Languages*. MIT Press, 2002.
- [14] PIERCE, B. C., AND TURNER, D. N. Local Type Inference. *ACM TOPLAS* (2000).
- [15] THIEMANN, P. J. Unboxed Values and Polymorphic Typing Revisited. In *Functional Programming Languages and Computer Architecture* (1995), ACM.
- [16] URECHE, V., BURMAKO, E., AND ODERSKY, M. Unifying Data Representation Transformations. Tech. rep., EPFL, 2014.
- [17] URECHE, V., TALAU, C., AND ODERSKY, M. Miniboxing: Improving the Speed to Code Size Tradeoff in Parametric Polymorphism Translations. In *OOPSLA* (2013).