

Provingly Correct Optimizations in an Intrinsically Typed Compiler

Experimentation Project Report

Matthias Heinzl (1632256)

August 3, 2022

Contents

1	Introduction	2
2	A Simple Expression Language	2
2.1	Expressions with Bindings	2
2.2	Representing Variables	3
2.3	Intrinsically Typed Syntax Trees	3
3	Dead Binding Elimination	5
3.1	Subsets of a Context	5
3.2	Live Variable Analysis	6
3.3	Transformation	7
3.4	Correctness	8
3.5	Iterating the Analysis	8
4	Further Work	9
4.1	Extending the Language	9
4.2	Other Analyses	9

1 Introduction

When writing a compiler for a programming language, an important consideration is the treatment of binders and variables. A well-known technique when using dependently typed programming languages such as Agda [6] is to define an intrinsically typed syntax tree, where expressions are scope- and type-safe by construction and admit a total evaluation function [2]. This construction has featured in several papers, exploring basic operations like renaming and substitution [1] as well as compilation to different target languages [7, supplemental material].

Performing optimisations on intrinsically typed programs, on the other hand, has not received as much attention. However, optimisations play an important role in compilers and establishing their correctness is often not trivial, with ample opportunity for mistakes. In this setting, program *analysis* not only needs to identify optimisation opportunities, but provide a proof witness that the optimisation is safe, e.g. that some dead code is indeed not used. For the *transformation* of the intrinsically typed program, the programmer can then rely on the compiler to check the relevant invariants, but it can be cumbersome to make it sufficiently clear that type- and scope-safety are preserved, especially when manipulating binders and variables.

As a step towards a more general treatment of optimisations of intrinsically typed programs, we present an implementation of *dead binding elimination* for a simple language. It first annotates expressions with variable usage information and then removes bindings that turn out to be unused. We further prove that the optimisation is semantics-preserving. The Agda source code is available online at <https://git.science.uu.nl/m.h.heinzel/correct-optimisations>.

2 A Simple Expression Language

2.1 Expressions with Bindings

We define a simple, typed expression language with let-bindings, variables, primitive values (integers and Booleans), and a few binary operators. Since the optimisations we are interested in relate to variables and binders only, the choice of possible values and additional primitive operations on them is mostly arbitrary. Extending the language with further values and operators is trivial.

$$\begin{array}{l} P, Q ::= v \\ \quad \mid P + Q \\ \quad \mid \text{let } x = P \text{ in } Q \\ \quad \mid x \end{array}$$

Expressions can be bound to a variable x using the **let** construction. Note that this makes the language equivalent to a restricted version of the simply typed λ -calculus, where λ -abstraction and application can only occur together as $(\lambda x.Q)P$. Encapsulating this pattern as **let** $x = P$ **in** Q simplifies parts of the analysis and avoids the need for allowing functions as values.

2.2 Representing Variables

The syntax specified above treats variables as letters, or more generally strings. To prevent complications with bindings of the same variable name shadowing each other and to make equality of terms independent of the specific names chosen (*α -equivalence*), compilers often represent variables in a different way. A popular choice are *de Bruijn indices*, where each variable is represented by a natural number, counting the number of bindings between variable occurrence and its binding: 0 refers to the innermost binding, 1 to the next-innermost etc.

Still, there might be *free variables*, where the de Bruijn index is larger than the number of bindings it has access to (*in scope*). If this happens unexpectedly during evaluation, an error is raised. Similarly, the type of a bound expression might not match the expected type at the variable occurrence where it is referenced. This makes the evaluation function partial; it should only be called after validating type- and scope-safety.

2.3 Intrinsically Typed Syntax Trees

When implementing a compiler in a dependently typed programming language, one does not need to accept partiality and the need for human vigilance. With *intrinsically typed syntax trees*, type- and scope-safety invariants are specified on the type level and verified by the type checker. We will demonstrate the approach in Agda and start by defining the types that terms can have.

```
data U : Set where
  BOOL : U
  NAT   : U

[ ] : U → Set
[ BOOL ] = Bool
[ NAT ] = Nat
```

To know if a variable occurrence is valid, we must consider its *context*, the bindings that are in scope. With de Bruijn indices in an untyped setting, it would suffice to know the number of bindings in scope. In a typed settings, we also need to know the type of each binding, so we represent the context by a list of types: One for each binding in scope, from innermost to outermost.

$\text{Ctx} = \text{List } \mathbf{U}$

variable
 $\Gamma : \text{Ctx}$
 $\sigma \ \tau : \mathbf{U}$

During evaluation, each variable in scope has a value. Together, we call these an *environment* in a given context.

data $\text{Env} : \text{Ctx} \rightarrow \text{Set}$ where
 $\text{Nil} : \text{Env } []$
 $\text{Cons} : [\![\sigma]\!] \rightarrow \text{Env } \Gamma \rightarrow \text{Env } (\sigma :: \Gamma)$

A variable then is an index into its context, also guaranteeing that its type matches that of the binding. Since variable $\text{Ref } \sigma \ \Gamma$ acts as a proof that the environment $\text{Env } \Gamma$ contains an element of type σ , variable lookup is total.

data $\text{Ref } (\sigma : \mathbf{U}) : \text{Ctx} \rightarrow \text{Set}$ where
 $\text{Top} : \text{Ref } \sigma \ (\sigma :: \Gamma)$
 $\text{Pop} : \text{Ref } \sigma \ \Gamma \rightarrow \text{Ref } \sigma \ (\tau :: \Gamma)$
 $\text{lookup} : \text{Ref } \sigma \ \Gamma \rightarrow \text{Env } \Gamma \rightarrow [\![\sigma]\!]$
 $\text{lookup } \text{Top} \ (\text{Cons } v \ \text{env}) = v$
 $\text{lookup } (\text{Pop } i) \ (\text{Cons } v \ \text{env}) = \text{lookup } i \ \text{env}$

Now we can define intrinsically typed expressions, where an Expr is indexed by both its type $(\sigma : \mathbf{U})$ and context $(\Gamma : \text{Ctx})$. We can see how the context changes when introducing a new binding that is then available in the body of a Let .

data $\text{Expr } (\Gamma : \text{Ctx}) : (\sigma : \mathbf{U}) \rightarrow \text{Set}$ where
 $\text{Val} : [\![\sigma]\!] \rightarrow \text{Expr } \Gamma \ \sigma$
 $\text{Plus} : \text{Expr } \Gamma \ \text{NAT} \rightarrow \text{Expr } \Gamma \ \text{NAT} \rightarrow \text{Expr } \Gamma \ \text{NAT}$
 $\text{Let} : (\text{decl} : \text{Expr } \Gamma \ \sigma) \rightarrow (\text{body} : \text{Expr } (\sigma :: \Gamma) \ \tau) \rightarrow \text{Expr } \Gamma \ \tau$
 $\text{Var} : \text{Ref } \sigma \ \Gamma \rightarrow \text{Expr } \Gamma \ \sigma$

This allows the definition of a total evaluator using an environment matching the expression's context.

$\text{eval} : \text{Expr } \Gamma \ \sigma \rightarrow \text{Env } \Gamma \rightarrow [\![\sigma]\!]$
 $\text{eval } (\text{Val } v) \ \text{env} = v$
 $\text{eval } (\text{Plus } e_1 \ e_2) \ \text{env} = \text{eval } e_1 \ \text{env} + \text{eval } e_2 \ \text{env}$
 $\text{eval } (\text{Let } e_1 \ e_2) \ \text{env} = \text{eval } e_2 \ (\text{Cons } (\text{eval } e_1 \ \text{env}) \ \text{env})$
 $\text{eval } (\text{Var } x) \ \text{env} = \text{lookup } x \ \text{env}$

3 Dead Binding Elimination

Note that an expression is not forced to make use of the whole context to which it has access. Specifically, a let-binding introduces a new element into the context, but it might never be used in the body. One commonly wants to identify such unused bindings so they can be removed from the program [5, Live Variable Analysis]. To that end, we annotate terms with the subset of their context that they actually use.

3.1 Subsets of a Context

A *sub-context* is not just any context. All its elements must come from the original context and we also require that their order is unchanged. This notion is known as a *thinning* or *order-preserving embedding* (OPE) [4] between the two contexts. We combine the sub-context with its OPE in a single data type `Subset`, which for each element of the context specifies whether to keep it or not. The uniquely determined sub-context can be extracted using `[_]`.

```
data Subset : Ctx → Set where
  Empty : Subset []
  Drop   : Subset Γ → Subset (τ :: Γ)
  Keep   : Subset Γ → Subset (τ :: Γ)

variable
  Δ Δ' Δ₁ Δ₂ : Subset Γ

[_] : Subset Γ → Ctx
[ Empty ]           = []
[ Drop Δ ]          = [ Δ ]
[ Keep {Γ} {τ} Δ ] = τ :: [ Δ ]
```

We can then define a subset-like order \subseteq :

```
_⊆_ : Subset Γ → Subset Γ → Set
Δ₁ ⊆ Keep Δ₂ = pop Δ₁ ⊆ Δ₂
Empty ⊆ Empty = ⊤
Drop Δ₁ ⊆ Drop Δ₂ = Δ₁ ⊆ Δ₂
Keep Δ₁ ⊆ Drop Δ₂ = ⊥
```

Dealing with `Subset` will simplify our task, especially the correctness proofs. While an OPE itself could be used to express a similar relation between contexts, its inhabitants are not unique. For example, `[NAT]` can be embedded into `[NAT , NAT]` either by keeping the first element or the second. This would complicate equality proofs on terms including witnesses of the relation. For similar reasons, it is unclear how we would want to define a union operator `_∪_ : Ctx → Ctx → Ctx`, but it is trivial on `Subsets`.

From now on, we will only consider expressions `Expr [Δ] σ` in some sub-context. Initially, we take $\Delta = \text{all } \Gamma : \text{Subset } \Gamma$, the complete sub-context of the original context.

3.2 Live Variable Analysis

The next step is to annotate each expression with its *live variables*, the sub-context Δ' that is really used. To that end, we define annotated expressions `LiveExpr $\Delta \Delta' \sigma$` . Note that Δ is treated as Γ before, conceptually accumulating all defined variables top-down. On the other hand, Δ' accumulates used variables bottom-up, starting with a singleton sub-context at the variable occurrences.

```
data LiveExpr { $\Gamma$  : Ctx} : ( $\Delta \Delta'$  : Subset  $\Gamma$ )  $\rightarrow$  ( $\sigma$  : U)  $\rightarrow$  Set where
  Val :
    [  $\sigma$  ]  $\rightarrow$ 
    LiveExpr  $\Delta \emptyset \sigma$ 
  Plus :
    LiveExpr  $\Delta \Delta_1$  NAT  $\rightarrow$ 
    LiveExpr  $\Delta \Delta_2$  NAT  $\rightarrow$ 
    LiveExpr  $\Delta (\Delta_1 \cup \Delta_2)$  NAT
  Let :
    LiveExpr  $\Delta \Delta_1 \sigma \rightarrow$ 
    LiveExpr { $\sigma :: \Gamma$ } (Keep  $\Delta$ )  $\Delta_2 \tau \rightarrow$ 
    LiveExpr  $\Delta (\Delta_1 \cup \text{pop } \Delta_2) \tau$ 
  Var :
    ( $x$  : Ref  $\sigma$  [  $\Delta$  ])  $\rightarrow$ 
    LiveExpr  $\Delta (\text{sing } \Delta x) \sigma$ 
```

To create such annotated expressions, we need to perform some static analysis of our source programs. The function `analyse` computes the live sub-context Δ' together with a matching annotated expression.

```
-- decide which variables are used or not
analyse :  $\forall \Delta \rightarrow$  Expr [  $\Delta$  ]  $\sigma \rightarrow \Sigma [ \Delta' \in \text{Subset } \Gamma ] \text{LiveExpr } \Delta \Delta' \sigma$ 
analyse  $\Delta$  (Val  $v$ ) =  $\emptyset$  , Val  $v$ 
analyse  $\Delta$  (Plus  $e_1 e_2$ ) with analyse  $\Delta e_1$  | analyse  $\Delta e_2$ 
... |  $\Delta_1$  ,  $le_1$  |  $\Delta_2$  ,  $le_2 = (\Delta_1 \cup \Delta_2)$  , Plus  $le_1 le_2$ 
analyse  $\Delta$  (Let  $e_1 e_2$ ) with analyse  $\Delta e_1$  | analyse (Keep  $\Delta$ )  $e_2$ 
... |  $\Delta_1$  ,  $le_1$  |  $\Delta_2$  ,  $le_2 = (\Delta_1 \cup \text{pop } \Delta_2)$  , Let  $le_1 le_2$ 
analyse  $\Delta$  (Var  $x$ ) = sing  $\Delta x$  , Var  $x$ 
```

The only requirement we have for it is that we can forget the annotations again.

```
-- forget the information about variable usage
forget : LiveExpr  $\Delta \Delta' \sigma \rightarrow$  Expr [  $\Delta$  ]  $\sigma$ 

-- forget  $\circ$  analyse  $\equiv$  id
analyse-preserves : ( $e$  : Expr [  $\Delta$  ]  $\sigma$ )  $\rightarrow$  forget (proj2 (analyse  $\Delta e$ ))  $\equiv e$ 
```

3.3 Transformation

It is useful to realise that we can evaluate `LiveExpr` directly, almost identically to the previous `eval` function. The main difference is in the `Let`-case, where we match on Δ_2 to distinguish whether the bound variable is live. If it is not, we can directly evaluate the body, ignoring the bound declaration.

Another important detail is that evaluation works under an environment for any sub-context Δ_u containing (at least) the live variables Δ' . Requiring an `Env [Δ]` would force us to provide a value even for unused variables. Fixing it to exactly `Env [Δ']`, however, would make it necessary to project the environment for each recursive call, since the subexpressions generally only use parts of the context.

```

evalLive :  $\forall \Delta_u \rightarrow \text{LiveExpr } \Delta \Delta' \tau \rightarrow \text{Env } [ \Delta_u ] \rightarrow .(\Delta' \subseteq \Delta_u) \rightarrow [ \tau ]$ 
evalLive  $\Delta_u$  (Val  $v$ ) env  $H = v$ 
evalLive  $\Delta_u$  (Plus  $\{\Delta\} \{\Delta_1\} \{\Delta_2\} e_1 e_2$ ) env  $H =$ 
  evalLive  $\Delta_u$   $e_1$  env ( $\subseteq_{U_1}$ -trans  $\Delta_1 \Delta_2 \Delta_u H$ )
  + evalLive  $\Delta_u$   $e_2$  env ( $\subseteq_{U_2}$ -trans  $\Delta_1 \Delta_2 \Delta_u H$ )
evalLive  $\Delta_u$  (Let  $\{\Delta = \Delta\} \{\Delta_1 = \Delta_1\} \{\Delta_2 = \text{Drop } \Delta_2\} e_1 e_2$ ) env  $H =$ 
  evalLive (Drop  $\Delta_u$ )  $e_2$  env ( $\subseteq_{U_2}$ -trans  $\Delta_1 \Delta_2 \Delta_u H$ )
evalLive  $\Delta_u$  (Let  $\{\Delta = \Delta\} \{\Delta_1 = \Delta_1\} \{\Delta_2 = \text{Keep } \Delta_2\} e_1 e_2$ ) env  $H =$ 
  evalLive (Keep  $\Delta_u$ )  $e_2$ 
  (Cons (evalLive  $\Delta_u$   $e_1$  env ( $\subseteq_{U_1}$ -trans  $\Delta_1 \Delta_2 \Delta_u H$ )) env)
  ( $\subseteq_{U_2}$ -trans  $\Delta_1 \Delta_2 \Delta_u H$ )
evalLive  $\{\Gamma\} \{\Delta\} \Delta_u$  (Var  $x$ ) env  $H = \text{lookupLive } \Delta \Delta_u x \text{ env } H$ 

```

This *optimised semantics* hints at a similar program transformation and will be useful when proving it to be correct. The transformation simply maps each constructor to its counterpart in `Expr`, with some variable renaming (e.g. from `[Δ_1]` to `[$\Delta_1 \cup \Delta_2$]`) and the abovementioned case distinction, where we can omit unused bindings.

```

dbe : LiveExpr  $\Delta \Delta' \sigma \rightarrow \text{Expr } [ \Delta' ] \sigma$ 
dbe (Val  $v$ ) =
  Val  $v$ 
dbe (Plus  $\{\Delta\} \{\Delta_1\} \{\Delta_2\} e_1 e_2$ ) =
  Plus (injExpr1  $\Delta_1 \Delta_2$  (dbe  $e_1$ )) (injExpr2  $\Delta_1 \Delta_2$  (dbe  $e_2$ ))
dbe (Let  $\{\Delta_1 = \Delta_1\} \{\Delta_2 = \text{Drop } \Delta_2\} e_1 e_2$ ) =
  injExpr2  $\Delta_1 \Delta_2$  (dbe  $e_2$ )
dbe (Let  $\{\Delta_1 = \Delta_1\} \{\Delta_2 = \text{Keep } \Delta_2\} e_1 e_2$ ) =
  Let
    (injExpr1  $\Delta_1 \Delta_2$  (dbe  $e_1$ ))
    (renameExpr (Keep  $\Delta_2$ ) (Keep ( $\Delta_1 \cup \Delta_2$ )) ( $\subseteq_{U_2}$   $\Delta_1 \Delta_2$ ) (dbe  $e_2$ ))
dbe  $\{\Gamma\} \{\Delta\} (\text{Var } x) =$ 
  Var (sing-ref  $\Delta x$ )

```

```

sing-ref : ( $\Delta : \text{Subset } \Gamma$ ) ( $x : \text{Ref } \sigma [ \Delta ]$ )  $\rightarrow \text{Ref } \sigma [ \text{sing } \Delta x ]$ 

```

As opposed to the type signature of `forget`, which returns to the original context, in `dbe` we drop unused variables, only keeping `[Δ']`.

3.4 Correctness

We want to show that dead binding elimination preserves semantics: $\text{eval} \circ \text{dbe} \circ \text{analyse} \equiv \text{eval}$. Since we know that $\text{forget} \circ \text{analyse} \equiv \text{id}$, it is sufficient to show $\text{eval} \circ \text{dbe} \equiv \text{eval} \circ \text{forget}$. The proof gets simpler if we split it up using the optimised semantics:

$$\text{eval} \circ \text{dbe} \equiv \text{evalLive} \equiv \text{eval} \circ \text{forget}$$

The actual proof statements are more involved, since they quantify over the expression and environment used. As foreshadowed in the definition of `evalLive`, the statements are also generalised to evaluation under any $\text{Env} \lfloor \Delta_u \rfloor$, as long as it contains the live sub-context. This gives us more flexibility when using the inductive hypothesis.

```
-- eval ∘ dbe ≡ evalLive
dbe-correct :
  (e : LiveExpr Δ Δ' σ) (Δu : Subset Γ) (env : Env ⌊ Δu ⌋) →
  .(H : Δ' ⊆ Δu) →
  eval (renameExpr Δ' Δu H (dbe e)) env ≡ evalLive Δu e env H

-- evalLive ≡ eval ∘ forget
evalLive-correct :
  (e : LiveExpr Δ Δ' σ) (Δu : Subset Γ) (env : Env ⌊ Δ ⌋) →
  .(H' : Δ' ⊆ Δu) → .(H : Δu ⊆ Δ) →
  evalLive Δu e (prjEnv Δu Δ H env) H' ≡ eval (forget e) env
```

Both proofs work inductively on the expression, with most cases being a straight-forward congruence. The interesting one is again `Let`, where we split cases on the variable being used or not and need some auxiliary facts about evaluation, renaming and sub-contexts.

An important detail is that the \subseteq proofs are *irrelevant arguments*, indicated by the preceding dot. This is possible since their witnesses are unique and makes it clear to the Agda type checker that e.g. $\text{prjEnv } \Delta' \Delta H_1$ and $\text{prjEnv } \Delta' \Delta H_2$ are equal.

3.5 Iterating the Analysis

The declaration of a binding that is removed can contain the only occurrences of some other variable. This makes another binding dead, allowing further optimisation when running the algorithm again. While in our simple setting all these bindings could be identified in a single pass using *strong live variable analysis*, in general it can be useful to simply iterate the optimisation until a fixpoint is reached.

This approach is not structurally recursive and the number of iterations is not known upfront, so Agda's termination checker needs our help. We

observe that the algorithm must terminate, since the number of bindings decreases with each iteration (except for the last) and cannot become negative. This corresponds to the ascending chain condition in program analysis literature [5]. To convince the termination checker, we use *well-founded recursion* [3] on the number of bindings.

The overall correctness follows directly from the correctness of each individual iteration step.

4 Further Work

4.1 Extending the Language

Since our language only contains let-bindings, it might be of interest to extend it with λ -abstractions (forming a simply-typed λ -calculus). Some increase in complexity seems necessary to eliminate applications of functions that do not use their argument, but we hope that our work is still largely applicable. The problem gets more challenging when introducing recursive bindings. Conversely, adding sum and product types might require more extensive bookkeeping, but should not pose fundamental difficulties.

4.2 Other Analyses

There are several other binding-related transformations to explore, such as moving bindings up or down in the syntax tree. Another interesting type of optimisation is avoidance of redundant computations using *available expression analysis*. An example is *common subexpression elimination*, where subexpressions get replaced by variables bound to equivalent declarations (pre-existing or newly created).

Ideally, further exploration will lead to the discovery of common patterns and useful strategies for performing optimisations on intrinsically typed syntax trees.

References

- [1] Guillaume Allais, Robert Atkey, James Chapman, Conor McBride, and James McKinna. A type and scope safe universe of syntaxes with binding: Their semantics and proofs. *Proc. ACM Program. Lang.*, 2(ICFP), jul 2018.
- [2] Lennart Augustsson and Magnus Carlsson. An exercise in dependent types: A well-typed interpreter. In *Workshop on Dependent Types in Programming*, 1999.
- [3] Ana Bove, Alexander Krauss, and Matthieu Sozeau. Partiality and recursion in interactive theorem provers – an overview. *Mathematical Structures in Computer Science*, 26(1):38–88, 2016.
- [4] James Maitland Chapman. *Type checking and normalisation*. PhD thesis, University of Nottingham, 2009.
- [5] Flemming Nielson, Hanne Riis Nielson, and Chris Hankin. *Principles of Program Analysis*. Springer Berlin, Heidelberg, Germany, 1st edition, 2014.
- [6] Ulf Norell. *Towards a practical programming language based on dependent type theory*. PhD thesis, Chalmers University of Technology, 2007.
- [7] Mitchell Pickard and Graham Hutton. Calculating dependently-typed compilers (functional pearl). *Proc. ACM Program. Lang.*, 5(ICFP), aug 2021.