

Type Preserving Compilation in Haskell

Stefan Monnier

DIRO - Université de Montréal

Motivation: why show type-preservation?

Representing CPS for STLC in Haskell

Showing CPS's type preservation

Extending to System F

Alternative representations

Program verification

Many techniques exist to *verify* properties exhaustively

- Hoare logic: Encode in logic the semantics of the language, the specification of the program, then prove correctness
- Automatic program derivation: Write a clean and naive version of the program, plus automatically-verified transformations
- Model checking: Abstract a model of some aspect of the code; exhaustively and automatically check it against a specification
- Ad-hoc analyses: Look for known error patterns; often imprecise
- Type checking: *modular* verification against a predefined set of rules

Usually can't be all used at the same time

Type checking

By far the most popular formal method

Fully integrated in the toolchain: “no extra work”

Good error reporting, directly pointing to offending line

Early detection \Rightarrow cheap fix

Guaranteed 100% pure formal: checks the *actual* code

Example, to eliminate NULL dereferences:

```
data Maybe a = Nothing | Just a
```

- `String`: type of strings; can't be “NULL”!
- `Maybe String`: this one can be “NULL”!

Trusted Computing Base (TCB)

Type checking can verify important properties about *source code*

But is the machine code really faithful to your source code?

What about the runtime functions?

Usually all these are part of the TCB

The compiler and runtime are complex programs, hard to trust

Proving their full semantic correctness is hard and costly

Type preservation

Rather than ensure full correctness, concentrate on *type preservation*

- Input and output terms of code transformations have related type

Model-checking or Hoare-style proofs might still be broken

But type-based properties of source code are faithfully preserved!

Sweet spot:

- Ensuring type preservation should be cheap(er)
- Most programs have no formal property to preserve other than types

Type preservation

Rather than ensure full correctness, concentrate on *type preservation*

- Input and output terms of code transformations have related type

Model-checking or Hoare-style proofs might still be broken

But type-based properties of source code are faithfully preserved!

Sweet spot:

- Ensuring type preservation should be cheap(er)
- Most programs have no formal property to preserve other than types

Note: do not confuse the object compiler and the meta compiler!

Typed Intermediate Languages

Don't throw away type info after type checking

Make all intermediate languages typed

Perform type-checking after each compilation phase

- Extra cost to store and compute the type information
- Extra cost of running the type checker
- Lost optimization opportunities
- Amounts to testing type-preservation
- Can only point to a whole compilation phase as the culprit

Type-Based Verification

Type checking can eliminate run-time errors, of course

It can also provide more guarantees about behavior:

- `iton: Int → String` guarantees we will not get NULL
- `id: α → α` guarantees preservation of type

Rich type systems can enforce precise invariants about programs

Proofs are intertwined with the program

When the proof's structure matches the program's

⇒ most of the proof (the boring part) comes for free

Verifying type preservation

Express type preservation in the types. I.e., instead of:

$$e : SrcExp \longrightarrow \boxed{\text{Code Transformation}} \longrightarrow e' : DstExp$$

Use:

$$e : SrcExp \tau \longrightarrow \boxed{\text{Code Transformation}} \longrightarrow e' : DstExp \tau$$

$e : SrcExp \tau$ means:

- e is a Haskell expression that evaluates to an AST
- That AST represents a source expression of type τ

We can use type-checking to verify the type preservation!

Type-based verification vs TIL

Expected advantages:

- No extra time spent manipulating types while running the compiler
- Errors detected earlier (during type-checking rather than at runtime)
- Errors linked directly to offending source code
- Potentially *less work* if proof mirrors code

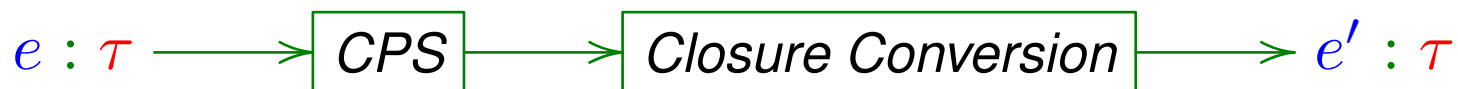
Might still cause lost optimization opportunities

Requires advanced type systems

Sometimes imposes a change in representation

Code transformations

Two common code transformations in compilers for functional languages



CPS:

- Only use tail-calls, equivalent to *goto*
- Make the return addresses and stack explicit

Closure conversion:

- Make closure construction explicit
- All functions are closed and can be hoisted to top-level

Why CPS and closure conversion

CPS is the ideal case:

- The proof mirrors the code exactly
- Code is never moved from one scope to another

Closure conversion is challenging:

- Safety depends on correctness of some parts
- Scoping is affected, so need to reason explicitly about environments
- Perceived to be impossible with HOAS

Will focus on CPS today

Source language

Source language is simply typed lambda calculus (STLC), e.g.:

$$\tau ::= \text{Int} \mid \text{Bool} \mid \tau_1 \rightarrow \tau_2$$

$$e ::= n \mid \text{if } e_1 e_2 e_3 \mid e_1 e_2 \mid \lambda x \rightarrow e \mid x$$

Could be encoded in Haskell as:

```
data Exp where
```

```
  Num  :: Int  -> Exp
```

```
  If    :: Exp -> Exp -> Exp -> Exp
```

```
  App   :: Exp -> Exp -> Exp
```

```
  Lam   :: Var -> Exp -> Exp
```

```
  Vref  :: Var -> Exp
```

Annotated source language

Annotate Exp with the source-level type of the expression

data $Exp \tau$ where

$Num :: Int \rightarrow Exp Int$

$If :: Exp Bool \rightarrow Exp \tau \rightarrow Exp \tau \rightarrow Exp \tau$

$App :: Exp (\tau_1 \rightarrow \tau_2) \rightarrow Exp \tau_1 \rightarrow Exp \tau_2$

... ..

We reuse/abuse Haskell types in red to denote *source* types

Does it really express a valid typing?

STLC typing rules

Judgment $\boxed{\vdash e : \tau}$ means e has type τ

$$\frac{}{\vdash n : \mathit{Int}} \quad \frac{\vdash e_1 : \mathit{Bool} \quad \vdash e_2 : \tau \quad \vdash e_3 : \tau}{\vdash \text{if } e_1 \ e_2 \ e_3 : \tau}$$

$$\frac{\vdash e_1 : \tau_1 \rightarrow \tau_2 \quad \vdash e_2 : \tau_1}{\vdash e_1 \ e_2 : \tau_2}$$

We can see that

$$e :: \mathit{Exp} \ \tau \quad \simeq \quad \vdash \llbracket e \rrbracket : \tau$$

Where $\llbracket e \rrbracket$ converts e from Haskell back to its STLC source

Typing rules in Haskell

$$\boxed{\vdash \llbracket e \rrbracket : \tau}$$

$$\boxed{e :: \text{Exp } \tau}$$

$$\frac{}{\vdash n : \text{Int}}$$

$$\text{Num} :: \text{Int} \rightarrow \text{Exp Int}$$

$$\frac{\vdash e_1 : \text{Bool} \quad \vdash e_2 : \tau \quad \vdash e_3 : \tau}{\vdash \text{if } e_1 \ e_2 \ e_3 : \tau}$$

$$\text{If} :: \text{Exp Bool} \rightarrow \text{Exp } \tau \rightarrow \text{Exp } \tau$$

$$\frac{\vdash e_1 : \tau_1 \rightarrow \tau_2 \quad \vdash e_2 : \tau_1}{\vdash e_1 \ e_2 : \tau_2}$$

$$\text{App} :: \text{Exp } (\tau_1 \rightarrow \tau_2) \rightarrow \text{Exp } \tau_1 \rightarrow \text{Exp } \tau_2$$

Curry-Howard

Curry-Howard isomorphism states:

$$\begin{array}{l} \text{Types} \quad \Leftrightarrow \quad \text{Propositions} \\ \text{Expressions} \quad \Leftrightarrow \quad \text{Proofs} \end{array}$$

In our case:

$$e :: \text{Exp } \tau \quad \simeq \quad \vdash \llbracket e \rrbracket : \tau$$

Our expressions e of type $\text{Exp } \tau$ play double role:

- Encode syntax of source term
- Encode typing derivation

We say they are *intrinsically typed*

Important simplification!

Encoding bindings

Compilers often represent variables as strings or unique integers

For practical typing we need something else

$$\text{Lam} :: \text{Var } \tau_1 \rightarrow \text{Exp } \tau_2 \rightarrow \text{Exp } (\tau_1 \rightarrow \tau_2)$$

$$\text{Vref} :: \text{Var } \tau \rightarrow \text{Exp } \tau$$

How do we make sure that a $\text{Var } \tau$ indeed refers to a var of type τ ?

- Higher-order abstract syntax: elegant when it works
- DeBruijn indices: frustrating but pragmatic and universal

Higher-order abstract syntax (HOAS)

Use meta-level variables to represent object-level variables:

$$\begin{array}{c}
 \vdash x : \tau_1 \\
 \vdots \\
 \vdash e : \tau_2 \\
 \hline
 \vdash \lambda x \rightarrow e : \tau_1 \rightarrow \tau_2
 \end{array}
 \qquad
 \begin{array}{l}
 \text{Lam} :: (\text{Exp } \tau_1 \rightarrow \text{Exp } \tau_2) \\
 \rightarrow \text{Exp } (\tau_1 \rightarrow \tau_2)
 \end{array}$$

No *Vref* case!

E.g. a source code like $\lambda f \rightarrow \lambda x \rightarrow f x$ would be represented as

$$\text{Lam } (\lambda f \rightarrow \text{Lam } (\lambda x \rightarrow \text{App } f x))$$

CPS transformation

Make the return addresses and stack explicit; Eliminate non-tail calls

Functions never return

$inc\ x = x + 1$

...

let $z = inc\ a$

in ...

becomes

$inc\ (x, k) = k\ (x + 1)$

...

$inc\ (a,$

$\lambda z \rightarrow \dots)$

Name all intermediate values

Explicit ordering of evaluation

Remaining (tail) calls correspond to simple branch instructions

Can be used to compile exception handling and call/cc

Sketching type preservation of CPS

$$\llbracket \tau_1 \rightarrow \tau_2 \rrbracket \Rightarrow (\llbracket \tau_1 \rrbracket, \llbracket \tau_2 \rrbracket \rightarrow \emptyset) \rightarrow \emptyset$$

$$\llbracket \text{Int} \rrbracket \Rightarrow \text{Int}$$

$$\llbracket x \rrbracket k \Rightarrow k x$$

$$\llbracket \lambda x \rightarrow e \rrbracket k \Rightarrow k (\lambda(x, k') \rightarrow \llbracket e \rrbracket k')$$

$$\llbracket e_1 e_2 \rrbracket k \Rightarrow \llbracket e_1 \rrbracket (\lambda x_1 \rightarrow \llbracket e_2 \rrbracket (\lambda x_2 \rightarrow x_1(x_2, k)))$$

Type preservation lemma looks like

$$\vdash e : \tau \Rightarrow \vdash \llbracket e \rrbracket : (\llbracket \tau \rrbracket \rightarrow \emptyset) \rightarrow \emptyset$$

Representing CPS converted code in Haskell

data *ValK* τ where

NumK :: *Int* \rightarrow *ValK Int*

PairK :: *ValK* $\tau_1 \rightarrow$ *ValK* $\tau_2 \rightarrow$ *ValK* (τ_1, τ_2)

LamK :: (*ValK* $\tau_1 \rightarrow$ *ExpK*) \rightarrow *ValK* $(\tau_1 \rightarrow \emptyset)$

... ..

data *ExpK* where

IfK :: *ValK Bool* \rightarrow *ExpK* \rightarrow *ExpK* \rightarrow *ExpK*

AppK :: *ValK* $(\tau_1 \rightarrow \emptyset) \rightarrow$ *ValK* $\tau_1 \rightarrow$ *ExpK*

... ..

CPS type preservation in Haskell

type family *CPS* τ — Implements $\llbracket \tau \rrbracket$

type instance *CPS* *Int* = *Int*

type instance *CPS* ($\tau_1 \rightarrow \tau_2$) = (*CPS* τ_1 , (*CPS* τ_2) $\rightarrow \emptyset$) $\rightarrow \emptyset$

cps :: *Exp* $\tau \rightarrow ((\text{ValK} ((\text{CPS } \tau) \rightarrow \emptyset)) \rightarrow \text{ExpK})$

Intrinsic typing:

- *Exp*, *ExpK*, *ValK* encode typing derivations (and syntax trees)
- *cps* encodes the proof of type preservation (and the transformation)

cps is an executable proof!

Sample of CPS transformation

$$\llbracket e_1 e_2 \rrbracket k \Rightarrow \llbracket e_1 \rrbracket (\lambda x_1 \rightarrow \llbracket e_2 \rrbracket (\lambda x_2 \rightarrow x_1(x_2, k)))$$

Ignoring a bit of gymnastic needed to work with HOAS in Haskell:

$$cps :: Exp \tau \rightarrow ((ValK ((CPS \tau) \rightarrow \emptyset)) \rightarrow ExpK)$$

$$cps (App e_1 e_2) k =$$

$$cps e_1 (\lambda v_1 \rightarrow$$

$$cps e_2 (\lambda v_2 \rightarrow$$

$$AppK v_1 (PairK v_2 (LamK k)))$$

Look ma! No types! Type preservation for free!

STLC is very limited \Rightarrow extend with parametric polymorphism:

$$\tau ::= \dots \mid \tau_1 \rightarrow \tau_2 \mid \forall t. \tau \mid t$$

$$e ::= \dots \mid e_1 e_2 \mid \lambda x \rightarrow e \mid x \mid \Lambda t. e \mid e[\tau]$$

2 new bindings and 1 new kind of variable

Using HOAS, it would be neatly encoded in Haskell as:

$$TLam :: (\text{forall } t. \text{Exp } \tau) \rightarrow \text{Exp } (\forall (\lambda t. \tau))$$

$$TApp :: \text{Exp } (\forall f) \rightarrow \text{Exp } (f \tau)$$

But GHC does not support type-level λ !

System F in Haskell

Use de Bruijn indices to represent the variables bound by \forall :

`data Z; data S τ` — To represent natural numbers

`data Var τ ; data All τ`

So $\forall t_1. \forall t_2. t_1 \rightarrow t_2$ is represented as `All (All ((S Z) \rightarrow Z))`

`TLam :: (forall t . Exp (Subst τ t Z)) \rightarrow Exp (All τ)`

`TApp :: Exp (All τ_1) \rightarrow Exp (Subst τ_1 τ_2 Z)`

We still use HOAS for the Λ binding

We have to implement the type-level function `Subst`

Substitution

de Bruijn indices use type-level naturals Z and $S \tau$

Subst requires type-level operation on naturals

In the *TApp* and *TLam* cases of *cps*, we need to show commutativity

$$CPS (Subst \tau_1 \tau_2 n) \equiv Subst (CPS \tau_1) (CPS \tau_2) n$$

GHC does not offer any way to prove it at the type-level:

subst_cps_commute :: *TypeRep* τ_1 \rightarrow *TypeRep* τ_2

\rightarrow *Equiv* (*CPS* (*Subst* $\tau_1 \tau_2 n$)) (*Subst* (*CPS* τ_1) (*CPS* τ_2) n)

If we could use HOAS, we'd get *Subst* and commutativity for free

CPS without CPS

CPS type-level function can be eliminated:

- Keep $\tau_1 \rightarrow \tau_2$ as the representation for “normal functions”
- Replace $\tau \rightarrow \emptyset$ with new $Cont \tau$ for continuations

We lose the equivalence $\tau_1 \rightarrow \tau_2 \equiv Cont (\tau_1, Cont \tau_2)$

This equivalence is not really used

We need two different $LamK$ and two different $AppK$

A bit of code duplication, but avoid need for $subst_cps_commute$

Still need to define $Subst$ and type-level naturals

CPS with de Bruijn

HOAS can't always be used (e.g. for closure conversion)

Other typing judgment $\boxed{\Gamma \vdash e : \tau}$

$$\frac{\Gamma, x:\tau_1 \vdash e : \tau_2}{\Gamma \vdash \lambda x \rightarrow e : \tau_1 \rightarrow \tau_2} \qquad \frac{\Gamma(x) = \tau}{\Gamma \vdash x : \tau}$$

Corresponding Haskell judgment $\boxed{e :: \text{Exp } \Gamma \ \tau}$

data *Exp* $\Gamma \ \tau$ where

Lam :: *Exp* (Γ, τ_1) $\tau_2 \rightarrow$ *Exp* Γ ($\tau_1 \rightarrow \tau_2$)

Vref :: *VarIn* $\Gamma \ \tau \rightarrow$ *Exp* $\Gamma \ \tau$

de Bruijn indices

Represent the proof of membership by the index in the set

```
data VarIn  $\Gamma$   $\tau$  where
  VarZ :: VarIn ( $\Gamma$ ,  $\tau$ )  $\tau$ 
  VarS :: VarIn  $\Gamma$   $\tau_1$   $\rightarrow$  VarIn ( $\Gamma$ ,  $\tau_2$ )  $\tau_1$ 
```

But Γ in source and Γ after CPS are different

Need extra rewriting and corresponding proofs

A lot more work than with HOAS

Had to use it for closure conversion

Conclusion

Type-preserving CPS and closure conversion of System F in Haskell

CPS part worked rather well

Closure conversion is painful

The code is executable version of the proof

Need better type system and more automation to make it practical