

Compiler Correctness via Contextual Equivalence

Matthew McKay, advised by Karl Cray

Abstract

We have developed a methodology for verifying the correctness of the closure conversion phase of a compiler, adapted from the work by Perconti and Ahmed. This lets us verify that individual components of programs are compiled correctly, so they can be linked with any other code and still behave as desired. We do this by using a combined language that encompasses both the source and target languages in which the compiled code can be reasoned about alongside its source, which we do using contextual equivalence. Our main improvement over previous methods is that we do not need boundaries that separate the source and target language while inside the combined language.

Combined Language

Our language C is a combined language of the source and target languages, with all terms and types from both languages. The source language S is System F extended with recursive functions and existentials. The target T is the closure converted version of the source, so it is the same except it has closed recursive functions and application instead of normal recursive functions and application.

$$\begin{aligned} \tau &::= \alpha \mid \text{unit} \mid \text{int} \mid \tau \times \tau \mid \tau \rightarrow \tau \mid \tau \Rightarrow \tau \mid \forall \alpha. \tau \mid \exists \alpha. \tau \\ e &::= () \mid n \mid e p e \mid \text{ifz}(e, e, e) \mid x \mid \langle e, e \rangle \mid \pi_i e \\ &\quad \mid \text{fun } f(x : \tau). e \mid e e \mid \overline{\text{fun}} f(x : \tau). e \mid e \wedge e \\ &\quad \mid \Lambda \alpha. e \mid e[\tau] \mid \text{pack}[\tau', e] \text{ as } \exists \alpha. \tau \mid \text{unpack}[\alpha, x] = e \text{ in } e \\ p &::= + \mid - \mid * \\ \Gamma &::= \cdot \mid \Gamma, x : \tau \\ \Delta &::= \cdot \mid \Delta, \alpha \end{aligned}$$

Two Different Functions

The main thing of note about the combined language is the presence of two types of recursive functions. There are normal recursive functions, and then there are closed recursive functions, which are generated from normal recursive functions during closure conversion. They behave the same as normal recursive functions, however they can only use their argument and bound function variable, all other variables are not in scope inside the function (hence "closed").

$$\frac{\Delta \vdash \tau \text{ type} \quad \Delta; \Gamma, f : \tau \rightarrow \tau', x : \tau \vdash e : \tau'}{\Delta; \Gamma \vdash \text{fun } f(x : \tau). e : \tau \rightarrow \tau'} Tfun$$

$$\frac{\Delta \vdash \tau \text{ type} \quad \Delta; f : \tau \rightarrow \tau', x : \tau \vdash e : \tau'}{\Delta; \Gamma \vdash \overline{\text{fun}} f(x : \tau). e : \tau \Rightarrow \tau'} Tccfun$$

Closure Conversion

The closure conversion translation is essentially the standard translation, with the slight difference that the functions in the translated language are denoted as closed recursive functions. Thus we have the type translation for functions as

$$|\tau_1 \rightarrow \tau_2| = \exists \alpha. ((|\tau_1| \times \alpha) \Rightarrow |\tau_2|) \times \alpha$$

For closure conversion the only interesting translation rules are for functions and applications, the rules for which are below:

$$\frac{\Gamma = x_1 : \tau_1, \dots, x_n : \tau_n \quad \Delta \vdash_S \tau \text{ type} \quad \Delta; \Gamma, x : \tau, f : \tau \rightarrow \tau' \vdash_S e : \tau' \rightsquigarrow \bar{e} \quad \tau_{env} = |\tau_1| \times \dots \times |\tau_n|}{\Delta; \Gamma \vdash_S \text{fun } f(x : \tau). e : \tau \rightarrow \tau' \rightsquigarrow \text{pack}[\tau_{env}, ((\overline{\text{fun}} f(y : |\tau| \times \tau_{env}). [F/f]S(\bar{e})), E)] \text{ as } |\tau \rightarrow \tau'|} Rfun$$

$$\frac{\Delta; \Gamma \vdash_S e_1 : \tau \rightarrow \tau' \rightsquigarrow \bar{e}_1 \quad \Delta; \Gamma \vdash_S e_2 : \tau \rightsquigarrow \bar{e}_2}{\Delta; \Gamma \vdash_S e_1 e_2 : \tau' \rightsquigarrow \text{unpack}[\alpha, x] = \bar{e}_1 \text{ in } (\pi_1 x) \wedge \langle \bar{e}_2, \pi_2 x \rangle} Rapp$$

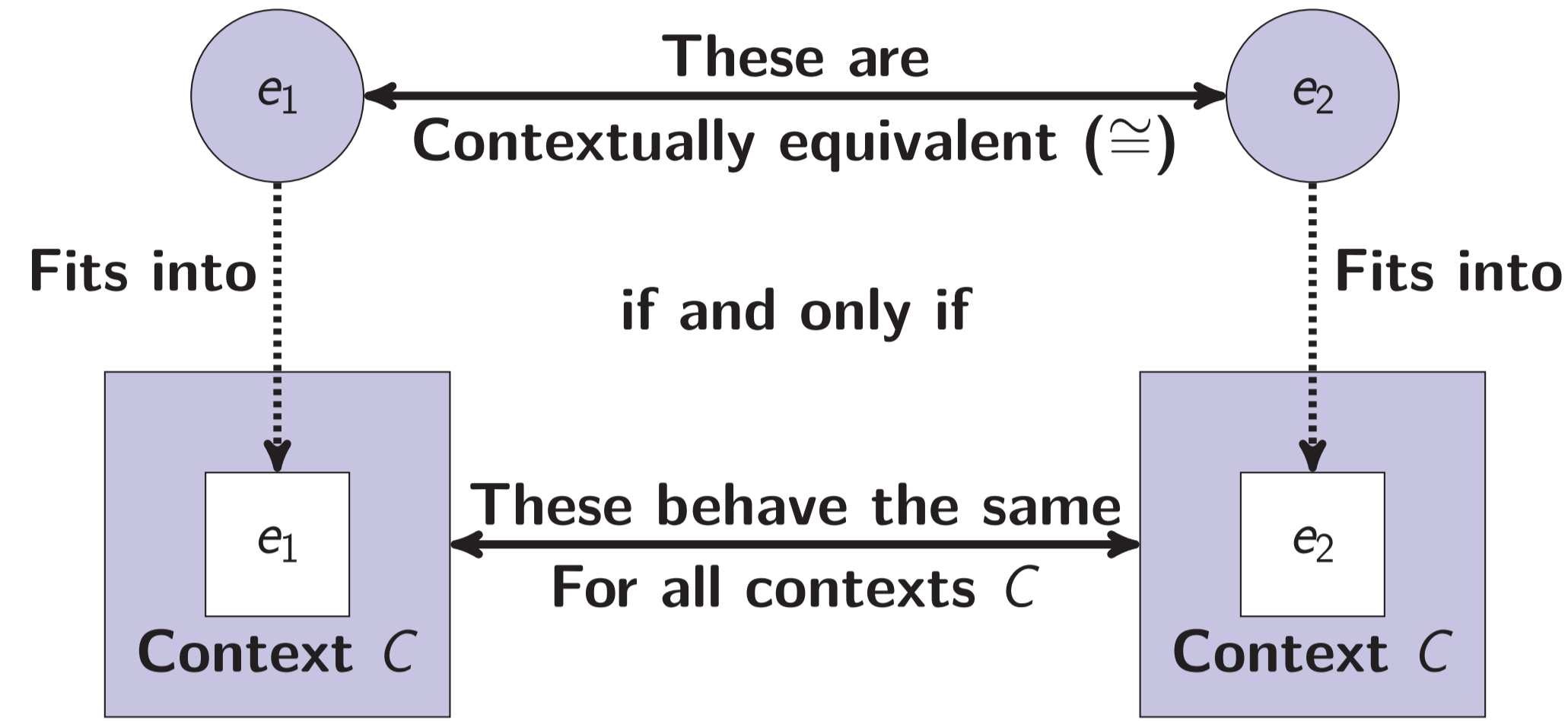
The first rule uses the following definitions:

$$\begin{aligned} E &= \langle x_1, \langle \dots \langle x_{n-1}, x_n \dots \rangle \dots \rangle \rangle \\ S &= [\pi_1 y / x] [\pi_1 \pi_2 y / x_1] \dots [\pi_1 \pi_2 \dots \pi_2 y / x_{n-1}] [\pi_2 \dots \pi_2 y / x_n] \\ F &= \text{pack}[\tau_{env}, \langle f, \pi_2 y \rangle] \text{ as } |\tau \rightarrow \tau'| \end{aligned}$$

Contextual Equivalence

The idea behind contextual equivalence is that for two terms to be equivalent, any program that could use them will behave the same no matter which term is used inside it. Formally, we write $\Delta; \Gamma \vdash_C e_1 \cong e_2 : \tau$ if both terms e_1 and e_2 can be typed at τ , and for every context $C : (\Delta; \Gamma \triangleright \tau) \rightsquigarrow (\cdot \triangleright \text{int})$ (a program with a hole in it which both e_1 and e_2 fit into), we have that $\mathcal{C}\{e_1\} \simeq \mathcal{C}\{e_2\}$, which simply states that $\mathcal{C}\{e_1\}$ and $\mathcal{C}\{e_2\}$ will have the same terminating behavior (either loop or terminate).

Contextual Equivalence Visualization



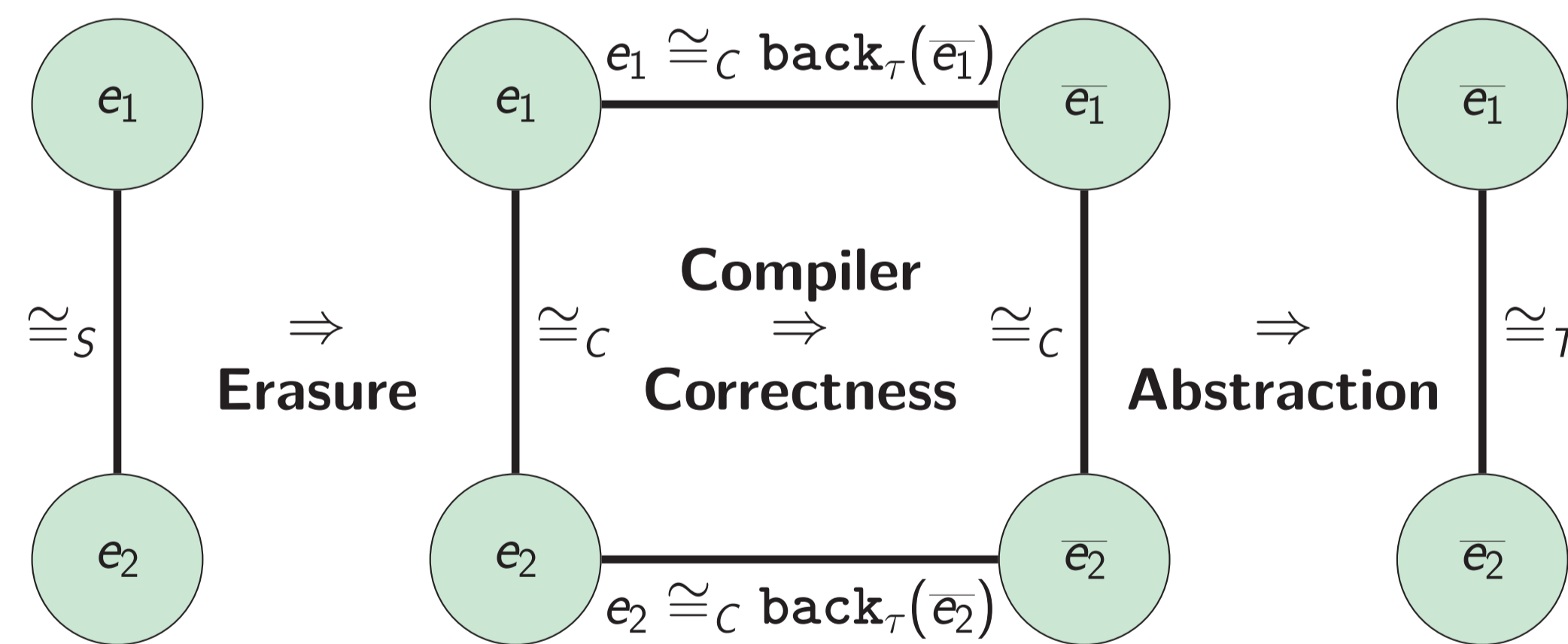
Statement of Compiler Correctness Theorem

Given a term $\Delta; \Gamma \vdash_S e_1 : \tau$ whose translation is $\Delta; \Gamma \vdash_S e_1 : \tau \rightsquigarrow \bar{e}_1$, we have that

$$\Delta; \Gamma \vdash e_1 \vdash \text{back}_\tau([\text{over}_\Gamma / \Gamma] \bar{e}_1) : \tau$$

This essentially says that e_1 is contextually equivalent to \bar{e}_1 , however these two terms don't have the same type in the combined language. Thus we need to wrap the compiled \bar{e}_1 in an over_Γ substitution and back_τ to make it have a matching type. Thus this theorem guarantees that the compiler preserves the behavior of the code that it compiles, so we have a verification of its correctness.

Equivalence Preservation Visualization



The over_τ substitutions for compiler correctness are left out for conciseness.

Statement of Equivalence Preservation Theorem

Given terms $\Delta; \Gamma \vdash_S e_1 : \tau$ and $\Delta; \Gamma \vdash_S e_2 : \tau$ whose translations are

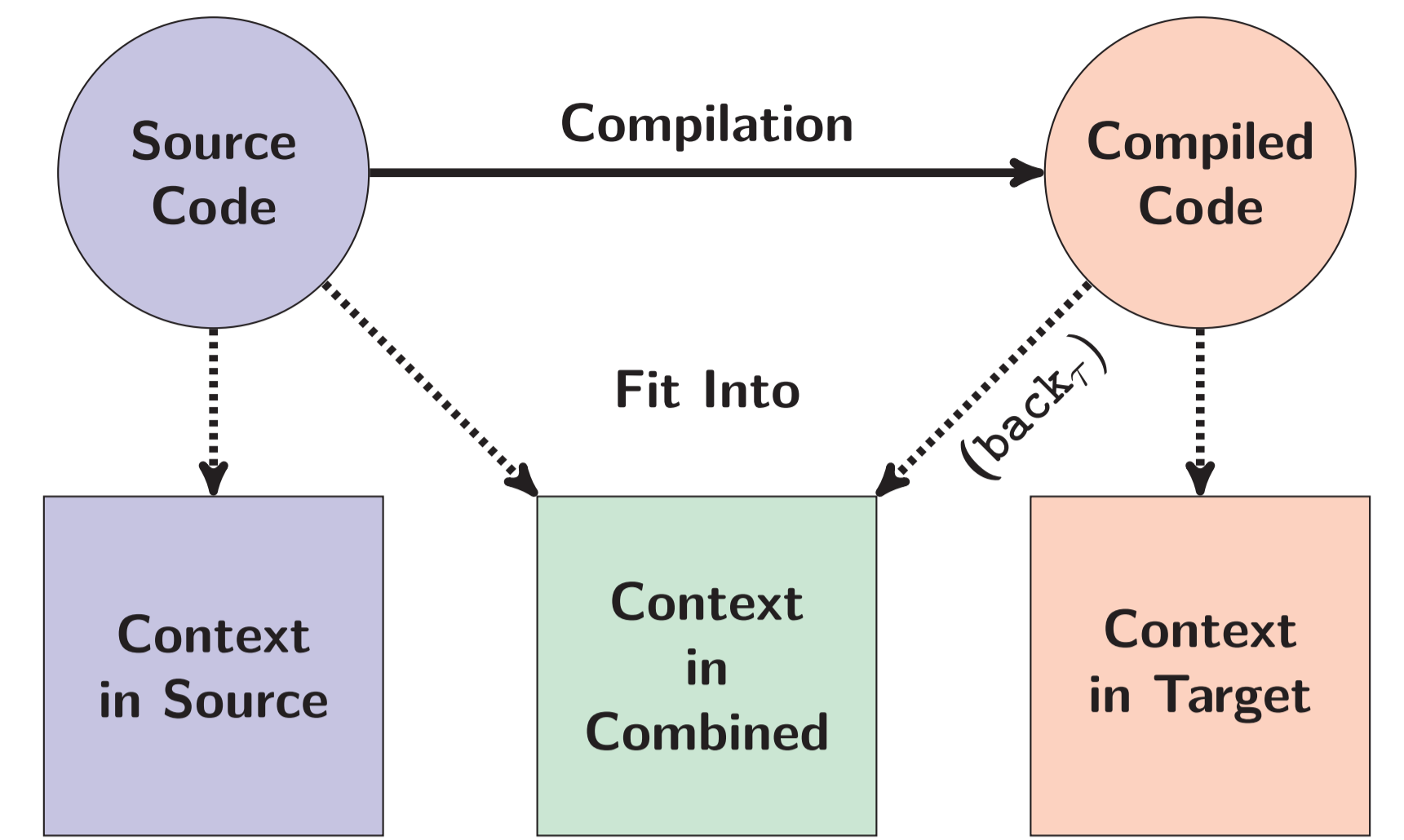
$\Delta; \Gamma \vdash_S e_1 : \tau \rightsquigarrow \bar{e}_1$ and $\Delta; \Gamma \vdash_S e_2 : \tau \rightsquigarrow \bar{e}_2$, then we have that

$$\Delta; \Gamma \vdash_S e_1 \cong e_2 : \tau \Rightarrow \Delta; \Gamma \vdash_T \bar{e}_1 \cong \bar{e}_2 : \tau$$

This essentially states that if two terms are equivalent in the source language, then the terms that they compile to are equivalent in the target language. Therefore, the compiler preserves equivalence.

To prove this we use our compiler correctness theorem, but we also need to show that if two terms are equivalent in S then they are equivalent in C , and similarly if two terms are equivalent in C then they are equivalent in T . The latter is simple since T is a subset of C , but for the former we need another translation, which we call erasure.

Combined Language Context Visualization



Over and Back Functions

We define two functions in the combined language, $\text{over}_\tau : \tau \rightarrow |\tau|$ and $\text{back}_\tau : |\tau| \rightarrow \tau$, each of which take terms in either the source or the target language and changes them at the top level to have the correct converted type in the opposite language. The two functions are mutually recursive and are inverses of one another.

$$\begin{aligned} \text{over}_\alpha &= \lambda x : \alpha. x \\ \text{over}_{\text{unit}} &= \lambda x : \text{unit}. x \\ \text{over}_{\text{int}} &= \lambda x : \text{int}. x \\ \text{over}_{\tau_1 \times \tau_2} &= \lambda x : \tau_1 \times \tau_2. \langle \text{over}_{\tau_1} \pi_1 x, \text{over}_{\tau_2} \pi_2 x \rangle \\ \text{over}_{\tau_1 \rightarrow \tau_2} &= \lambda f : \tau_1 \rightarrow \tau_2. \text{pack}[\tau_1 \rightarrow \tau_2, \langle \lambda y : |\tau_1| \times (\tau_1 \rightarrow \tau_2). \\ &\quad \text{over}_{\tau_2}((\pi_2 y) (\text{back}_{\tau_1} \pi_1 y)), f \rangle] \text{ as } |\tau_1 \rightarrow \tau_2| \\ \text{over}_{\forall \alpha. \tau} &= \lambda x : (\forall \alpha. \tau). \Lambda \alpha. (\text{over}_\tau(x[\alpha])) \\ \text{over}_{\exists \alpha. \tau} &= \lambda x : (\exists \alpha. \tau). \text{unpack}[\alpha, y] = x \text{ in } (\text{pack}[\alpha, \text{over}_\tau(y)] \text{ as } |\exists \alpha. \tau|) \\ \text{back}_\alpha &= \lambda x : \alpha. x \\ \text{back}_{\text{unit}} &= \lambda x : \text{unit}. x \\ \text{back}_{\text{int}} &= \lambda x : \text{int}. x \\ \text{back}_{\tau_1 \times \tau_2} &= \lambda x : |\tau_1 \times \tau_2|. \langle \text{back}_{\tau_1} \pi_1 x, \text{back}_{\tau_2} \pi_2 x \rangle \\ \text{back}_{\tau_1 \rightarrow \tau_2} &= \lambda f : |\tau_1 \rightarrow \tau_2|. \lambda y : \tau_1. \text{unpack}[\alpha, g] = f \text{ in } \\ &\quad \text{back}_{\tau_2}((\pi_1 g) \wedge \langle \text{over}_{\tau_1} y, \pi_2 g \rangle) \\ \text{back}_{\forall \alpha. \tau} &= \lambda x : |\forall \alpha. \tau|. \Lambda \alpha. (\text{back}_\tau(x[\alpha])) \\ \text{back}_{\exists \alpha. \tau} &= \lambda x : |\exists \alpha. \tau|. \text{unpack}[\alpha, y] = x \text{ in } (\text{pack}[\alpha, \text{back}_\tau(y)] \text{ as } |\exists \alpha. \tau|) \end{aligned}$$

Erasure

The erasure translation e° simply removes closed recursive functions by converting them to normal recursive functions. So $(\tau_1 \rightarrow \tau_2)^\circ = \tau_1^\circ \rightarrow \tau_2^\circ$ and

$$(\overline{\text{fun}} f(x : \tau). e)^\circ = \text{fun } f(x : \tau^\circ). e^\circ$$

The point of erasure is to say that adding closed recursive functions to S doesn't give it any more capability than it already has, which makes sense as closed functions do exactly the same thing as normal functions. This can then be used to show that equivalent terms in S are equivalent in C , since removing the closed functions doesn't affect the language's power.

Future Work

There is plenty more work that can be extended from this research, as it was mostly a demonstration of a possible approach. The most logical extension would be to apply this methodology to other translation phases of a compiler, like the allocation phase. Unfortunately other phases will be significantly more complicated to do, as the combined language increases in complexity as the difference between the source and target languages increases. Closure conversion is a good starting place due to the similarity between the languages, though it could be possible to extend this methodology to other phases of a compiler.