

Space-Efficient Gradual Typing

David Herman¹, Aaron Tomb², and Cormac Flanagan²

¹ Northeastern University

² University of California, Santa Cruz

Abstract

Gradual type systems offer a smooth continuum between static and dynamic typing by permitting the free mixture of typed and untyped code. The runtime systems for these languages—and other languages with hybrid type checking—typically enforce function types by dynamically generating function proxies. This approach can result in unbounded growth in the number of proxies, however, which drastically impacts space efficiency and destroys tail recursion.

We present an implementation strategy for gradual typing that is based on *coercions* instead of function proxies, and which combines adjacent coercions to limit their space consumption. We prove bounds on the space consumed by coercions as well as soundness of the type system, demonstrating that programmers can safely mix typing disciplines without incurring unreasonable overheads. Our approach also detects certain errors earlier than prior work.

1 GRADUAL TYPING FOR SOFTWARE EVOLUTION

Dynamically typed languages have always excelled at exploratory programming. Languages such as Lisp, Scheme, Smalltalk, and JavaScript support quick early prototyping and incremental development without the overhead of documenting (often-changing) structural invariants as types. For large applications, however, static type systems have proven invaluable. They are crucial for understanding and enforcing key program invariants and abstractions, and they catch many errors early in the development process.

Given these different strengths, it is not uncommon to encounter the following scenario: a programmer builds a prototype in a dynamically-typed scripting language, perhaps even in parallel to a separate, official software development process with an entire team. The team effort gets mired in process issues and over-engineering, and the programmer’s prototype ends up getting used in production. Before long, this hastily conceived prototype grows into a full-fledged production system, but without the structure or guarantees provided by static types. The system becomes unwieldy; QA can’t produce test cases fast enough and bugs start cropping up that no one can track down. Ultimately, the team decides to port the application to a statically typed language, requiring a complete rewrite of the entire system.

The cost of cross-language migration is huge and often insupportable. But the scenario above is avoidable. Several languages combine static and dynamic typing, among them Boo [7], Visual Basic.NET [19], Sage [16], and PLT Scheme [22]. This approach of *hybrid typing*, where types are enforced with a combination of

static and dynamic checks, has begun to gain attention in the research community [5, 16, 21, 22, 19, 2]. Recently, Siek and Taha [21, ?] coined the slogan *gradual typing* for this important application of hybrid typing: the ability to implement both partially-conceived prototypes and mature, production systems in the same programming language by gradually introducing type discipline. Gradual typing offers the possibility of continuous software evolution from prototype to product, thus avoiding the huge costs of language migration.

Our recent experience in the working group on the JavaScript [8] language specification provides a more concrete example. JavaScript is a dynamically-typed functional programming language that is widely used for scripting user interactions in web browsers, and is a key technology in Ajax applications [15]. The enormous popularity of the language is due in no small part to its low barrier to entry; anyone can write a JavaScript program by copying and pasting code from one web page to another. Its dynamic type system and fail-soft runtime semantics allow programmers to produce something that *seems* to work with a minimum of effort. The increasing complexity of modern web applications, however, has motivated the addition of a static type system. The working group’s intention¹ is not to abandon the dynamically-typed portion of the language—because of its usefulness and to retain backwards compatibility—but rather to allow typing disciplines to interact via a hybrid system that supports gradual typing.

1.1 The Cost of Gradual Typing

Gradually-typed languages support both statically-typed and dynamically-typed code, and include runtime checks (or type casts) at the boundaries between these two typing disciplines, to guarantee that dynamically-typed code cannot violate the invariants of statically-typed code. To illustrate this idea, consider the following code fragment, which passes an untyped variable x into a variable y of type `Int`:

```
let x = true in ... let y: Int = x in ...
```

During compilation, the type checker inserts a dynamic type cast $\langle \text{Int} \rangle$ to enforce the type invariant on y ; at run-time, this cast detects the attempted type violation:

```
let x = true in ... let y: Int = ( $\langle \text{Int} \rangle$  x) in ...
→* Error: “failed cast”
```

Unfortunately, even these simple, first-order type checks can result in unexpected costs, as in the following example, where a programmer has added some type annotations to a previously untyped program:

```
even          =  $\lambda n: \text{Int}. \text{if } (n = 0) \text{ then true else odd } (n - 1)$ 
odd: Int → Bool =  $\lambda n: \text{Int}. \text{if } (n = 0) \text{ then false else even } (n - 1)$ 
```

¹The JavaScript specification is a work in progress, but gradual typing is a key design goal [9].

This program seems innocuous, but suffers from a space leak. Since *even* is dynamically typed, the result of each call to *even* ($n - 1$) must be cast to `Bool`, resulting in unbounded growth in the control stack and destroying tail recursion.

Additional complications arise when first-class functions cross the boundaries between typing disciplines. In general, it is not possible to check if an untyped function satisfies a particular static type. A natural solution is to wrap the function in a proxy that, whenever it is applied, casts its argument and result values appropriately, ensuring that the function is only observed with its expected type. This proxy-based approach is used heavily in recent literature [11, 12, 16, 18, 21], but has serious consequences for space efficiency.

As a simple example, consider the following program in continuation-passing style, where both mutually recursive functions take a continuation argument k , but only one of these arguments is annotated with a precise type:

$$\begin{aligned} \textit{even} &= \lambda n:\text{Int}. \lambda k:(? \rightarrow ?). \quad \text{if } (n = 0) \text{ then } (k \text{ true}) \text{ else } \textit{odd} (n - 1) k \\ \textit{odd} &= \lambda n:\text{Int}. \lambda k:(\text{Bool} \rightarrow \text{Bool}). \text{ if } (n = 0) \text{ then } (k \text{ false}) \text{ else } \textit{even} (n - 1) k \end{aligned}$$

Here, the recursive calls to *odd* and *even* quietly cast the continuation argument k with higher-order casts $\langle \text{Bool} \rightarrow \text{Bool} \rangle$ and $\langle ? \rightarrow ? \rangle$, respectively. This means that the function argument k is wrapped in an additional function proxy at each recursive call!

The flexibility promised by gradual typing can only be achieved if programmers are free to decide how precisely to type various parts of a program. This flexibility is lost if adding type annotations can accidentally trigger asymptotic changes in its space consumption. In short, existing implementation techniques for gradual typing suffer from unacceptable space leaks.

1.2 Space-Efficient Gradual Typing

We present an implementation strategy for gradual typing that overcomes these problems. Our approach hinges on the simple insight that when proxies accumulate at runtime, they often contain redundant information. In the higher-order example above, the growing chain of function proxies contains only two distinct components, $\text{Bool} \rightarrow \text{Bool}$ and $? \rightarrow ?$, which could be merged to the simpler but equivalent $\text{Bool} \rightarrow \text{Bool}$ proxy.

Type casts behave like *error projections* [?], which are closed under composition. However, the syntax of casts does not always compose; for example, there is no cast c such that $\langle c \rangle e = \langle \text{Int} \rangle \langle \text{Bool} \rangle e$. Furthermore, projections are idempotent, which should allow us to eliminate duplicate casts. For example, $\langle \text{Bool} \rightarrow \text{Bool} \rangle \langle \text{Bool} \rightarrow \text{Bool} \rangle e = \langle \text{Bool} \rightarrow \text{Bool} \rangle e$. But such transformations are inconvenient, if not impossible, with a representation of higher-order type casts as functions.

Our formalization instead leverages Henglein’s *coercion calculus* [17], which provides a syntax for projections, called *coercions*, which are closed under a composition operator. This allows us to combine adjacent coercions in order to eliminate redundant information and thus guarantee clear bounds on space consumption. By

eagerly combining coercions, we can also detect certain errors immediately as soon as a function cast is applied; in contrast, prior approaches would not detect these errors until the casted function is invoked.

Our approach is applicable to many hybrid-typed languages [12, 13, 16] that use function proxies and hence are prone to space consumption problems. For clarity, we formalize our approach for the simply-typed λ -calculus with references. Of course, gradual typing is not restricted to such simple type systems: the Sage language [16] incorporates gradual typing as part of a very expressive type system with polymorphic functions, type operators, first-class types, dependent types, general refinement types, etc. Concurrently, Siek and Taha [21] developed gradual typing for the simpler language $\lambda_{\perp}^?$, which we use as the basis for this presentation.

The presentation of our results proceeds as follows. The following section reviews the syntax and type system of $\lambda_{\perp}^?$. Section 3 introduces the coercion algebra underlying our approach. Section 4 describes how we compile source programs into a target language with explicit coercions. Section 5 provides an operational semantics for that target language, and Section 6 proves bounds on the space consumption. Section 7 extends our approach to detect errors earlier and provide better coverage. The last two sections place our work into context.

2 GRADUALLY-TYPED LAMBDA CALCULUS

This section reviews the gradually-typed λ -calculus $\lambda_{\perp}^?$. This language is essentially the simply-typed λ -calculus extended with the type $?$ to represent dynamic types; it also includes mutable reference cells to demonstrate the gradual typing of assignments.

$$\begin{array}{l} \text{Terms:} \quad e ::= k \mid x \mid \lambda x:T. e \mid e e \mid \text{ref } e \mid !e \mid e := e \\ \text{Types:} \quad S, T ::= B \mid T \rightarrow T \mid ? \mid \text{Ref } T \end{array}$$

Terms include the usual constants, variables, abstractions, and applications, as well as reference allocation, dereference, and assignment. Types include the dynamic type $?$, function types $T \rightarrow T$, reference types $\text{Ref } T$, and some collection of ground or base types B (such as Int or Float).

The $\lambda_{\perp}^?$ type system is a little unusual in that it is based on an intransitive *consistency* relation $S \sim T$ instead of the more conventional, transitive subtyping relation $S <: T$. Any type is consistent with the type $?$, from which it follows that, for example, $\text{Bool} \sim ?$ and $? \sim \text{Int}$. However, booleans cannot be used directly as integers, which is why the consistency relation is not transitively closed. We do not assume the consistency relation is symmetric, since a language might, for example, allow coercions from integers to floats but not vice-versa.

The consistency relation is defined in Figure 1. Rules [C-DYNL] and [C-DYNR] allow all coercions to and from type $?$. The rule [C-FLOAT] serves as an example of asymmetry by allowing coercion from Int to Float but not the reverse. The rule [C-FUN] is reminiscent of the contravariant/covariant rule for function subtyping.

Figure 1: Source Language Type System

<u>Consistency rules</u>				$S \sim T$
[C-REFL]	[C-DYNR]	[C-DYNL]	[C-FLOAT]	
$\overline{T \sim T}$	$\overline{T \sim ?}$	$\overline{? \sim T}$	$\overline{\text{Int} \sim \text{Float}}$	
[C-FUN]		[C-REF]		
$\frac{T_1 \sim S_1 \quad S_2 \sim T_2}{(S_1 \rightarrow S_2) \sim (T_1 \rightarrow T_2)}$		$\frac{T \sim S \quad S \sim T}{\text{Ref } S \sim \text{Ref } T}$		
<u>Type rules</u>				$E \vdash e : T$
[T-VAR]	[T-FUN]	[T-CONST]		
$\frac{(x : T) \in E}{E \vdash x : T}$	$\frac{E, x : S \vdash e : T}{E \vdash (\lambda x : S. e) : (S \rightarrow T)}$	$\frac{}{E \vdash k : \text{ty}(k)}$		
[T-APP1]		[T-APP2]		
$\frac{E \vdash e_1 : (S \rightarrow T) \quad E \vdash e_2 : S' \quad S' \sim S}{E \vdash (e_1 e_2) : T}$		$\frac{E \vdash e_1 : ? \quad E \vdash e_2 : S}{E \vdash (e_1 e_2) : ?}$		
[T-REF]		[T-DEREF1]	[T-DEREF2]	
$\frac{E \vdash e : T}{E \vdash \text{ref } e : \text{Ref } T}$		$\frac{E \vdash e : \text{Ref } T}{E \vdash !e : T}$	$\frac{E \vdash e : ?}{E \vdash !e : ?}$	
[T-ASSIGN1]			[T-ASSIGN2]	
$\frac{E \vdash e_1 : \text{Ref } T \quad E \vdash e_2 : S \quad S \sim T}{E \vdash e_1 := e_2 : S}$			$\frac{E \vdash e_1 : ? \quad E \vdash e_2 : T}{E \vdash e_1 := e_2 : ?}$	

We extend the invariant reference cells of $\lambda_{\sim}^?$ to allow coercion from $\text{Ref } S$ to $\text{Ref } T$ via rule [C-REF], provided S and T are symmetrically consistent. Unlike functions, reference cells do not distinguish their output (“read”) type from their input (“write”) type, so coercion must be possible in either direction. For example, the two reference types Ref Int and $\text{Ref } ?$ are consistent.

Figure 1 also presents the type rules for the source language, which are mostly standard. Notice the presence of two separate rules for procedure application. Rule [T-APP1] handles the case where the operator is statically-typed as a function; in this case, the argument may have any type consistent with the function’s domain. Rule [T-APP2] handles the case where the operator is dynamically-typed, in which case the argument may be of any type. The two rules for assignment follow an analogous pattern, accepting a consistent type when the left-hand side is known to have type $\text{Ref } T$, and any type when the left-hand side is dynamically-typed. Similarly, dereference expressions only produce a known type when the argument has a reference type.

3 COERCIONS

To achieve a space-efficient implementation, we compile source programs into a target language with explicit coercions, drawn from Henglein’s theory of dynamic typing [17]. Like casts, coercions allow values of one type to be used at any consistent type, performing conversions if necessary. Their key benefit over prior proxy-based representations is that they are *combinable*; if two coercions are wrapped around a function value, then they can be safely combined into a single coercion, thus reducing the space consumption of the program without changing its semantic behavior.

The coercion language and its typing rules are both defined in Figure 2. The coercion judgment $\vdash c : S \rightsquigarrow T$ states that coercion c serves to coerce values from type S to type T . The identity coercion I (of type $\vdash c : T \rightsquigarrow T$ for any T) always succeeds. Conversely, the failure coercion `Fail` always fails. For each dynamic type tag D there is an associated tagging coercion $D!$ that produces values of type $?$, and a corresponding check-and-untag coercion $D?$ that takes values of type $?$. Thus, for example, we have $\vdash \text{Int}! : \text{Int} \rightsquigarrow ?$ and $\vdash \text{Int}? : ? \rightsquigarrow \text{Int}$.

The function checking coercion $\text{Fun}?$ converts a value of type $?$ to have the dynamic function type $? \rightarrow ?$. If a more precise function type is required, this value can be further coerced via a function coercion $\text{Fun } c \ d$, where c coerces function arguments and d coerces results. For example, the coercion $(\text{Fun } \text{Int}? \ \text{Int}!)$ coerces from $? \rightarrow ?$ to $\text{Int} \rightarrow \text{Int}$, by untagging function arguments (via $\text{Int}?$) and tagging function results (via $\text{Int}!$). Reference coercions also contain two components: the first for coercing values put into the reference cell; the second for coercing values read from the cell. Finally, the coercion $c;d$ represents coercion composition, *i.e.*, the coercion c followed by coercion d .

This coercion language is sufficient to translate between all consistent types: if types S and T are consistent, then the following partial function $\text{coerce}(S, T)$ is defined and returns the appropriate coercion between these types.

$$\begin{aligned}
 \text{coerce} : \text{Type} \times \text{Type} &\rightarrow \text{Coercion} \\
 \text{coerce}(T, T) &= I \\
 \text{coerce}(B, ?) &= B! \\
 \text{coerce}(?, B) &= B? \\
 \text{coerce}(\text{Int}, \text{Float}) &= \text{IntFloat} \\
 \text{coerce}(S_1 \rightarrow S_2, T_1 \rightarrow T_2) &= \text{Fun } \text{coerce}(T_1, S_1) \ \text{coerce}(S_2, T_2) \\
 \text{coerce}(?, T_1 \rightarrow T_2) &= \text{Fun}?; \text{coerce}(? \rightarrow ?, T_1 \rightarrow T_2) \\
 \text{coerce}(T_1 \rightarrow T_2, ?) &= \text{coerce}(T_1 \rightarrow T_2, ? \rightarrow ?); \text{Fun}! \\
 \text{coerce}(\text{Ref } S, \text{Ref } T) &= \text{Ref } \text{coerce}(T, S) \ \text{coerce}(S, T) \\
 \text{coerce}(?, \text{Ref } T) &= \text{Ref}?; \text{coerce}(\text{Ref } ?, \text{Ref } T) \\
 \text{coerce}(\text{Ref } T, ?) &= \text{coerce}(\text{Ref } T, \text{Ref } ?); \text{Ref}!
 \end{aligned}$$

Coercing a type T to itself produces the identity coercion I . Coercing base types B

Figure 2: Coercion Language and Type Rules

<i>Coercions:</i> $c, d ::= I \mid \text{Fail} \mid D! \mid D? \mid \text{IntFloat} \mid \text{Fun } c \ c \mid \text{Ref } c \ c \mid c; c$			
<i>Dynamic tags:</i> $D ::= B \mid \text{Fun} \mid \text{Ref}$			
<u>Coercion rules</u>			
<div style="border: 1px solid black; display: inline-block; padding: 2px 5px;">$\vdash c : S \rightsquigarrow T$</div>			
[C-ID]	[C-FAIL]	[C-B!]	[C-B?]
$\vdash I : T \rightsquigarrow T$	$\vdash \text{Fail} : S \rightsquigarrow T$	$\vdash B! : B \rightsquigarrow ?$	$\vdash B? : ? \rightsquigarrow B$
[C-FUN!]		[C-FUN?]	
$\vdash \text{Fun}! : (? \rightarrow ?) \rightsquigarrow ?$		$\vdash \text{Fun}? : ? \rightsquigarrow (? \rightarrow ?)$	
[C-FUN]			
$\vdash c_1 : T'_1 \rightsquigarrow T_1$		$\vdash c_2 : T_2 \rightsquigarrow T'_2$	
$\vdash (\text{Fun } c_1 \ c_2) : (T_1 \rightarrow T_2) \rightsquigarrow (T'_1 \rightarrow T'_2)$			
[C-REF!]		[C-REF?]	
$\vdash \text{Ref}! : (\text{Ref } ?) \rightsquigarrow ?$		$\vdash \text{Ref}? : ? \rightsquigarrow (\text{Ref } ?)$	
[C-REF]			
$\vdash c : T \rightsquigarrow S$		$\vdash d : S \rightsquigarrow T$	
$\vdash (\text{Ref } c \ d) : (\text{Ref } S) \rightsquigarrow (\text{Ref } T)$			
[C-COMPOSE]			[C-FLOAT]
$\vdash c_1 : T \rightsquigarrow T_1$		$\vdash c_2 : T_1 \rightsquigarrow T_2$	
$\vdash (c_1; c_2) : T \rightsquigarrow T_2$			$\vdash \text{IntFloat} : \text{Int} \rightsquigarrow \text{Float}$

to type $?$ requires a tagging coercion $B!$, and coercing $?$ to a base type B requires a runtime check $B?$. Function coercions work by coercing their domain and range types. The type $?$ is coerced to a function type via a two-step coercion: first the value is checked to be a function and then coerced from the dynamic function type $? \rightarrow ?$ to $T_1 \rightarrow T_2$. Dually, typed functions are coerced to type $?$ via coercion to a dynamic function type followed by the function tag $\text{Fun}!$. Coercing a $\text{Ref } S$ to a $\text{Ref } T$ entails coercing all writes from T to S and all reads from S to T . Coercing reference types to and from $?$ is analogous to function coercion.

Lemma 1 (Well-typed coercions).

1. $S \sim T$ iff $\text{coerce}(S, T)$ is defined.
2. If $c = \text{coerce}(S, T)$ then $\vdash c : S \rightsquigarrow T$.

Proof. Inductions on the derivations of $S \sim T$ and $\text{coerce}(S, T)$. \square

Figure 3: Target Language Syntax and Type Rules

<i>Terms:</i>	$s, t ::= x \mid f \mid k \mid a \mid \lambda x:T. t \mid t t \mid \text{ref } t \mid !t \mid t := t \mid \langle c \rangle t$
<i>Stores:</i>	$\sigma ::= \emptyset \mid \sigma[a := v]$
<i>Typing environments:</i>	$E ::= \emptyset \mid E, x : T$
<i>Store typings:</i>	$\Sigma ::= \emptyset \mid \Sigma, a : T$
<i>Values:</i>	$v ::= u \mid \langle c \rangle u$
<i>Uncoerced values:</i>	$u ::= f \mid k \mid a$
<i>Function closures:</i>	$f ::= \langle \lambda x:T. t, \rho \rangle$
<i>Environments:</i>	$\rho ::= \emptyset \mid \rho[x := v]$

Type rules	$E; \Sigma \vdash t : T \quad E; \Sigma \vdash \sigma$
------------	--

$\frac{[T\text{-VAR}] \quad (x:t) \in E}{E; \Sigma \vdash x : T}$	$\frac{[T\text{-FUN}] \quad E, x : S; \Sigma \vdash t : T}{E; \Sigma \vdash (\lambda x:S. t) : (S \rightarrow T)}$	$\frac{[T\text{-APP}] \quad E; \Sigma \vdash t_1 : (S \rightarrow T) \quad E; \Sigma \vdash t_2 : S}{E; \Sigma \vdash (t_1 t_2) : T}$
$\frac{[T\text{-REF}] \quad E; \Sigma \vdash t : T}{E; \Sigma \vdash \text{ref } t : \text{Ref } T}$	$\frac{[T\text{-DEREF}] \quad E; \Sigma \vdash t : \text{Ref } T}{E; \Sigma \vdash !t : T}$	$\frac{[T\text{-ASSIGN}] \quad E; \Sigma \vdash t_1 : \text{Ref } T \quad E; \Sigma \vdash t_2 : T}{E; \Sigma \vdash t_1 := t_2 : T}$
$\frac{[T\text{-CONST}] \quad}{E; \Sigma \vdash k : \text{ty}(k)}$	$\frac{[T\text{-CAST}] \quad \vdash c : S \rightsquigarrow T \quad E; \Sigma \vdash t : S}{E; \Sigma \vdash \langle c \rangle t : T}$	$\frac{[T\text{-ADDR}] \quad (a : T) \in \Sigma}{E; \Sigma \vdash a : \text{Ref } T}$
$\frac{[T\text{-CLOSURE}] \quad E; \Sigma \vdash \text{resolve}(f) : T}{E; \Sigma \vdash f : T}$	$\frac{[T\text{-STORE}] \quad \text{dom}(\sigma) = \text{dom}(\Sigma) \quad \forall a \in \text{dom}(\sigma). E; \Sigma \vdash \sigma(a) : \Sigma(a)}{E; \Sigma \vdash \sigma}$	

4 TARGET LANGUAGE AND COMPILATION

During compilation, we both type check the source program and insert explicit type casts where necessary. The target language of this compilation process is essentially the same as the source, except that it uses explicit coercions of the form $\langle c \rangle t$ as the only mechanism for connecting terms of type $?$ and terms of other types. For example, the term $\langle \text{Int?} \rangle x$ has type Int , provided that x has type $?$. The language syntax and type rules are defined in Figure 3, and are mostly straightforward.

The process of type checking and inserting coercions is formalized via the *compilation judgment*:

$$E \vdash e \hookrightarrow t : T$$

Here, the type environment E provides types for free variables, e is the original source program, t is a modified version of the original program with additional coercions, and T is the inferred type for t . The rules defining the compilation judgment are shown in Figure 4, and they rely on the partial function *coerce* to compute coercions between types. For example, rule [C-APP1] compiles an application expression

Figure 4: Compilation Rules

Compilation of terms	$E \vdash e \hookrightarrow t : T$	
$\frac{[C\text{-VAR}]}{(x : T) \in E} \frac{}{E \vdash x \hookrightarrow x : T}$	$\frac{[C\text{-CONST}]}{} \frac{}{E \vdash k \hookrightarrow k : \text{ty}(k)}$	$\frac{[C\text{-FUN}]}{} \frac{E, x : S \vdash e \hookrightarrow t : T}{E \vdash (\lambda x : S. e) \hookrightarrow (\lambda x : S. t) : (S \rightarrow T)}$
$\frac{[C\text{-APP1}]}{E \vdash e_1 \hookrightarrow t_1 : (S \rightarrow T) \quad E \vdash e_2 \hookrightarrow t_2 : S' \quad c = \text{coerce}(S', S)}{E \vdash e_1 e_2 \hookrightarrow (t_1 (\langle c \rangle t_2)) : T}$		
$\frac{[C\text{-APP2}]}{E \vdash e_1 \hookrightarrow t_1 : ? \quad E \vdash e_2 \hookrightarrow t_2 : S' \quad c = \text{coerce}(S', ?)}{E \vdash e_1 e_2 \hookrightarrow ((\langle \text{Fun?} \rangle t_1) (\langle c \rangle t_2)) : ?}$		
$\frac{[C\text{-REF}]}{E \vdash e \hookrightarrow t : T} \frac{}{E \vdash \text{ref } e \hookrightarrow \text{ref } t : \text{Ref } T}$	$\frac{[C\text{-DEREF1}]}{E \vdash e \hookrightarrow t : \text{Ref } T} \frac{}{E \vdash !e \hookrightarrow !t : T}$	$\frac{[C\text{-DEREF2}]}{E \vdash e \hookrightarrow t : ?} \frac{}{E \vdash !e \hookrightarrow !(\langle \text{Ref?} \rangle t) : ?}$
$\frac{[C\text{-ASSIGN1}]}{E \vdash e_1 \hookrightarrow t_1 : \text{Ref } S \quad E \vdash e_2 \hookrightarrow t_2 : T \quad c = \text{coerce}(T, S)}{E \vdash e_1 := e_2 \hookrightarrow (t_1 := (\langle c \rangle t_2)) : S}$		
$\frac{[C\text{-ASSIGN2}]}{E \vdash e_1 \hookrightarrow t_1 : ? \quad E \vdash e_2 \hookrightarrow t_2 : T \quad c = \text{coerce}(T, ?)}{E \vdash e_1 := e_2 \hookrightarrow ((\langle \text{Ref?} \rangle t_1) := (\langle c \rangle t_2)) : ?}$		

where the operator has a function type $S \rightarrow T$ by coercing the argument expression from type S' to S . Rule [C-APP2] handles the case where the operator is dynamically-typed by inserting a runtime check $\langle \text{Fun?} \rangle$ to ensure the operator evaluates to a function and coercing the argument to type $?$ to yield tagged values. Rules [C-REF] and [C-DEREF1] compile typed reference allocation and dereference. Rule [C-DEREF2] handles dynamically-typed dereference by inserting a runtime $\langle \text{Ref?} \rangle$ check. Rule [C-ASSIGN1] compiles statically-typed assignment, coercing the right-hand side to the expected type of the reference cell, and [C-ASSIGN2] compiles dynamically-typed assignment, coercing the left-hand side with a runtime $\langle \text{Ref?} \rangle$ check and the right-hand side with a runtime tagging coercion to type $?$. Compilation succeeds on all well-typed source programs, and produces only well-typed target programs.

Theorem 2 (Well-typed compilation). *For all E , e , and T , the following statements are equivalent:*

1. $E \vdash e : T$
2. $\exists t$ such that $E \vdash e \hookrightarrow t : T$ and $E; \emptyset \vdash t : T$

Proof. Inductions on the compilation and source language typing derivations. \square

5 OPERATIONAL SEMANTICS

We now consider how to implement the target language in a manner that limits the space consumed by coercions. The key idea is to combine adjacent coercions, thus eliminating redundant information while preserving the semantic behavior of programs.

Following Morrisett, Felleisen, and Harper [20], we use the CEKS machine [10] as an operational model of the target language that is representative of the space usage of realistic implementations. Figure 5 provides the definitions and rules of the abstract machine. Each configuration of the machine is either a term configuration $\langle t, \rho, C, \sigma \rangle$ or a value configuration $\langle v, C, \sigma \rangle$, representing evaluation of a term t in environment ρ or return of a value v , respectively, in an evaluation context C and store σ .

Most of the evaluation rules are standard for a context-machine reduction semantics. The most important rule is [E-CCAST], which ensures that adjacent casts are always merged:

$$\langle \langle c \rangle (\langle d \rangle t), \rho, C, \sigma \rangle \longrightarrow \langle \langle d; c \rangle t, \rho, C, \sigma \rangle \text{ [E-CCAST]}$$

In order to maintain bounds on their size, coercions are maintained normalized throughout evaluation according to the following rules:

$$\begin{array}{ll} I; c = c & D!; D? = I \\ c; I = c & D!; D'? = \text{Fail} \quad \text{if } D \neq D' \\ \text{Fail}; c = \text{Fail} & (\text{Fun } c_1 c_2); (\text{Fun } d_1 d_2) = \text{Fun } (d_1; c_1) (c_2; d_2) \\ c; \text{Fail} = \text{Fail} & (\text{Ref } c_1 c_2); (\text{Ref } d_1 d_2) = \text{Ref } (d_1; c_1) (c_2; d_2) \end{array}$$

This normalization is applied in a transitive, compatible manner whenever the rule [E-CCAST] is applied, thus bounding the size of coercions generated during evaluation. Evaluation satisfies the usual progress and preservation lemmas.

Theorem 3 (Soundness of evaluation). *If $\emptyset; \emptyset \vdash t : T$ then either*

1. $\langle t, \emptyset, \bullet, \emptyset \rangle$ diverges,
2. $\langle t, \emptyset, \bullet, \emptyset \rangle \longrightarrow^* \langle \langle \text{Fail} \rangle u, \rho, C, \sigma \rangle$, or
3. $\langle t, \emptyset, \bullet, \emptyset \rangle \longrightarrow^* \langle v, \bullet, \sigma \rangle$ and $\exists \Sigma$ such that $\emptyset; \Sigma \vdash v : T$ and $\emptyset; \Sigma \vdash \sigma$.

Proof. Via standard subject reduction and progress lemmas in the style of Wright and Felleisen [24]. \square

6 SPACE EFFICIENCY

We now consider how much space coercions consume at runtime, beginning with an analysis of how much space each individual coercion can consume.

Figure 5: Operational Semantics

<i>Configurations:</i>	$M ::= \langle t, \rho, C, \sigma \rangle \mid \langle v, C, \sigma \rangle$	
<i>Contexts:</i>	$C ::= \bullet \mid C \cdot \langle F, \rho \rangle$	
<i>Context Frames:</i>	$F ::= (\bullet t) \mid (v \bullet) \mid \text{ref } \bullet \mid !\bullet \mid \bullet := t \mid v := \bullet \mid \langle c \rangle \bullet$	
	$\langle x, \rho, C, \sigma \rangle \longrightarrow \langle \rho(x), C, \sigma \rangle$	[E-VAR]
	$\langle \lambda x:T. t, \rho, C, \sigma \rangle \longrightarrow \langle \langle \lambda x:T. t, \rho \rangle, C, \sigma \rangle$	[E-ABS]
	$\langle \langle c \rangle u, \rho, C, \sigma \rangle \longrightarrow \langle \langle c \rangle u, C, \sigma \rangle$ if $c \notin \{I, \text{Fail}\}$	[E-CAST]
	$\langle F[t], \rho, C, \sigma \rangle \longrightarrow \langle t, \rho, C \cdot \langle F, \rho \rangle, \sigma \rangle$ if $t \neq v$	[E-PUSH]
	$\langle v, C \cdot \langle F, \rho \rangle, \sigma \rangle \longrightarrow \langle F[v], \rho, C, \sigma \rangle$	[E-POP]
	$\langle \langle \lambda x:S. t, \rho \rangle v, _, C, \sigma \rangle \longrightarrow \langle t, \rho[x:=v], C, \sigma \rangle$	[E-BETA]
	$\langle k v, _, C, \sigma \rangle \longrightarrow \langle \delta(k, v), C, \sigma \rangle$	[E-PRIM]
	$\langle \text{ref } v, \rho, C, \sigma \rangle \longrightarrow \langle a, C, \sigma[a:=v] \rangle$ for $a \notin \text{dom}(\sigma)$	[E-NEW]
	$\langle !a, \rho, C, \sigma \rangle \longrightarrow \langle \sigma(a), C, \sigma \rangle$	[E-DEREF]
	$\langle a := v, \rho, C, \sigma \rangle \longrightarrow \langle v, C, \sigma[a:=v] \rangle$	[E-ASSIGN]
	$\langle \langle \langle \text{Fun } c d \rangle f \rangle v, \rho, C, \sigma \rangle \longrightarrow \langle \langle d \rangle (f (\langle c \rangle v)), \rho, C, \sigma \rangle$	[E-CBETA]
	$\langle !(\langle \text{Ref } c d \rangle a), \rho, C, \sigma \rangle \longrightarrow \langle \langle d \rangle !a, \rho, C, \sigma \rangle$	[E-CDEREF]
	$\langle \langle \langle \text{Ref } c d \rangle a \rangle := v, \rho, C, \sigma \rangle \longrightarrow \langle \langle d \rangle (a := \langle c \rangle v), \rho, C, \sigma \rangle$	[E-CASSIGN]
	$\langle \langle I \rangle u, \rho, C, \sigma \rangle \longrightarrow \langle u, C, \sigma \rangle$	[E-ID]
	$\langle \langle c \rangle (\langle d \rangle t), \rho, C, \sigma \rangle \longrightarrow \langle \langle d; c \rangle t, \rho, C, \sigma \rangle$	[E-CCAST]
	$\langle \langle \text{IntFloat} \rangle n, \rho, C, \sigma \rangle \longrightarrow \langle r, C, \sigma \rangle$ $r = \text{nearestFloat}(n)$	[E-FCAST]

6.1 Space Consumption

The size of a coercion $\text{size}(c)$ is defined as the size of its abstract syntax tree representation. When two coercions are sequentially composed and normalized during evaluation, the size of the normalized, composed coercion may of course be larger than either of the original coercions. In order to reason about the space required by such composed coercions, we introduce a notion of the *height* of a coercion:

$$\begin{aligned}
 \text{height}(I) &= \text{height}(\text{Fail}) = \text{height}(D!) = \text{height}(D?) = 1 \\
 \text{height}(\text{Ref } c d) &= \text{height}(\text{Fun } c d) = 1 + \max(\text{height}(c), \text{height}(d)) \\
 \text{height}(c; d) &= \max(\text{height}(c), \text{height}(d))
 \end{aligned}$$

Notably, the height of a composed coercion is bounded by the maximum height of its constituents. In addition, normalization never increases the height of a coercion. Thus, the height of any coercion created during program evaluation is never larger than the height of some coercion in the original compiled program.

Furthermore, this bound on the height of each coercion in turn guarantees a bound on the coercion's size, according to the following lemma. In particular, even though the length of a coercion sequence $c_1; \dots; c_n$ does not contribute to its height,

the restricted structure of well-typed, normalized coercions constrains the length (and hence size) of such sequences.

Lemma 4. *For all well-typed normalized coercions c , $size(c) \leq 5(2^{height(c)} + 1)$.*

Proof. Induction on c . Assume $c = c_1; \dots; c_n$, where each c_i is not a sequential composition. Suppose some $c_i = \text{Ref } d_1 \ d_2$. So $\vdash c_i : \text{Ref } S \rightsquigarrow \text{Ref } T$. Hence c_i can be preceded only by $\text{Ref}?$, which must be the first coercion in the sequence, and similarly can be followed only by $\text{Ref}!$, which must be the last coercion in the sequence.

Thus in the worst case $c = \text{Ref}?; \text{Ref } d_1 \ d_2; \text{Ref}!$ and $size(c) = 5 + size(d_1) + size(d_2)$. Applying the induction hypothesis to the sizes of d_1 and d_2 yields:

$$size(c) \leq 5 + 2(5(2^{height(c)-1} - 1)) = 5(2^{height(c)} - 1)$$

The case for $\text{Fun } d_1 \ d_2$ is similar. The coercions I and Fail can only appear alone. Finally, coercions of the form $D?; D!$ are valid. However, composition of a coercion c matching this pattern with one of the other valid coercions is either ill-typed or triggers a normalization that yields a coercion identical to c . \square

In addition, the height of any coercion created during compilation is bounded by the height of some type in the source program (where the height of a type is the height of its abstract syntax tree representation).

Lemma 5. *If $c = coerce(S, T)$, then $height(c) \leq \max(height(S), height(T))$.*

Proof (sketch). Induction on the structure of $coerce(S, T)$. \square

Theorem 6 (Size of coercions). *For any e , c such that*

1. $\emptyset \vdash e \hookrightarrow t : T$ and
2. $\langle t, \emptyset, \bullet, \emptyset \rangle \longrightarrow^* M$ and
3. M contains c ,

$\exists S \in e$ such that $size(c) \leq 5(2^{height(S)} - 1)$.

Proof. Induction on the length of the reduction sequence, using Lemma 4; the base case is by induction on the compilation derivation, using Lemma 5. \square

We now bound the total cost of maintaining coercions in the space-efficient semantics. We define the size of a configuration as the sum of the sizes of its components. In order to construct a realistic measure of the store, we count only those cells that an idealized garbage collector would consider live by restricting the $size$ function to the domain of reachable addresses:

$$\begin{aligned} size(\langle t, \rho, C, \sigma \rangle) &= size(t) + size(\rho) + size(C) + size(\sigma|_{reachable(t, \rho, C, \sigma)}) + 1 \\ size(\rho) &= \sum_{x \in dom(\rho)} (1 + size(\rho(x))) \\ size(\langle \lambda x : T. t, \rho \rangle) &= size(\lambda x : T. t) + size(\rho) + 1 \\ size(k) &= size(a) = 1 \\ &\dots \end{aligned}$$

To show that coercions occupy bounded space in the model, we compare the size of configurations in reduction sequences to configurations in an “oracle” semantics where coercions require no space. The oracular measure $size_{OR}$ is defined similarly to $size$, but without a cost for maintaining coercions:

$$\begin{aligned} size_{OR}(c) &= 0 \\ size_{OR}(\langle c \rangle t) &= size_{OR}(t) \\ size_{OR}(\langle c \rangle \bullet) &= 0 \\ &\dots \end{aligned}$$

The following theorem then bounds the fraction of the program state occupied by coercions in the space-efficient semantics.

Theorem 7 (Space consumption). *If $\emptyset \vdash e \hookrightarrow t : T$ and $\langle t, \emptyset, \bullet, \emptyset \rangle \longrightarrow^* M$, then there exists some $S \in e$ such that $size(M) \in O(2^{height(S)} \cdot size_{OR}(M))$.*

Proof. During evaluation, the reduction rules prevent the nesting of adjacent pairs of nested coercions in any value or term in the configuration; similarly, no evaluation context may contain adjacent nested coercions. Thus the number of coercions in any component of the configuration is proportional to the size of that component. By Theorem 6 the size of each coercion is in $O(2^{height(S)})$ for the largest S in e . \square

6.2 Tail recursion

Theorem 7 has the important consequence that coercions do not affect the control space consumption of tail-recursive programs. For example, the *even* and *odd* functions mentioned in the introduction now consume constant space. This important property is achieved by immediately combining adjacent coercions on the stack via the [E-CCAST] rule. This section sketches three implementation techniques for this rule.

Coercion-passing style This approach adds an extra argument to every procedure, representing the result coercion. Tail calls coalesce but do not perform this coercion, instead passing it along to the next function.

Trampoline A trampoline [14] is a well-known technique for implementing tail-recursive languages where tail calls are implemented by returning a thunk to a top-level loop. Tail-recursive functions with coercions return both a thunk and a coercion to the driver loop, which accumulates and coalesces returned coercions.

Continuation marks Continuation marks [6] allow programs to annotate continuation frames with arbitrary data. When a marked frame performs a tail call, the subsequent frame can inherit and modify the destroyed frame’s marks. Coercions on the stack are stored as marks and coalesced on tail calls.

7 EARLY ERROR DETECTION

Consider the following code fragment, which erroneously attempts to convert an $(\text{Int} \rightarrow \text{Int})$ function to have type $(\text{Bool} \rightarrow \text{Int})$:

```
let f :? = ( $\lambda x : \text{Int}. x$ ) in
let g :  $(\text{Bool} \rightarrow \text{Int}) = f$  in ...
```

Prior strategies for gradual typing would not detect this error until g is called.

In contrast, our coercion-based implementation allows us to detect this error as soon as g is defined. In particular, after compilation and evaluation, the value of g is

$$\langle \text{Fun Fail } I \rangle (\lambda x : \text{Int}. x)$$

where the domain coercion `Fail` explicates the inconsistency of the two domain types `Int` and `Bool`. We can modify our semantics to halt as soon as such inconsistencies are detected, by adding the following coercion normalization rules:

$$\begin{array}{ll} \text{Fun } c \text{ Fail} & = \text{Fail} & \text{Ref } c \text{ Fail} & = \text{Fail} \\ \text{Fun Fail } c & = \text{Fail} & \text{Ref Fail } c & = \text{Fail} \end{array}$$

Using these rules, our implementation strategy halts as soon as g is defined, resulting in earlier error detection and better coverage, since g may not actually be called in some tests.

8 RELATED WORK

There is a large body of work combining static and dynamic typing. The simplest approach is to use reflection with the type `?`, as in Amber [3]. Since case dispatch cannot be precisely type-checked with reflection alone, many languages provide statically-typed `typecase` on dynamically-typed values, including Simula-67 [1] and Modula-3 [4].

For dynamically-typed languages, *soft typing* systems provide type-like static analyses that facilitate optimization and early error reporting [23]. These systems may provide static type information but do not allow explicit type annotations, whereas enforcing documented program invariants (i.e., types) is a central feature of gradual typing.

Similarly, Henglein’s theory of dynamic typing [17] provides a framework for static type optimizations but only in a purely dynamically-typed setting. We use Henglein’s coercions instead for structuring the coercion algebra of our target language. Our application is essentially different: in the gradually-typed setting, coercions serve to enforce explicit type annotations, whereas in the dynamically-typed setting, coercions represent checks required by primitive operations.

None of these approaches facilitates *migration* between dynamically and statically-typed code, at best requiring hand-coded interfaces between them. The gradually-typed approach, exemplified by the gradually typed language $\lambda_{\rightarrow}^?$, of Siek and Taha [21],

lowers the barrier for code migration by allowing mixture of expressions of type τ with more precisely-typed expressions. Our work improves gradual typing by eliminating the drastic effects on space efficiency subtly incurred by crossing the boundary between typing disciplines.

Several other systems employ dynamic function proxies, including hybrid type checking [12], software contracts [11], and recent work on software migration by Tobin-Hochstadt and Felleisen [22]. We believe our approach to coalescing redundant proxies could improve the efficiency of all of these systems.

9 CONCLUSION AND FUTURE WORK

We have presented a space-efficient implementation strategy for the gradually-typed λ -calculus. More work remains to demonstrate the applicability of this technique in the setting of more advanced type systems. In particular, recursive types and polymorphic types may present a challenge for maintaining constant bounds on the size of coercions. We intend to explore techniques for representing these infinite structures as finite graphs.

Another useful feature for runtime checks is *blame annotations* [11], which pinpoint the particular expressions in the source program that cause coercion failures at runtime by associating coercions with the expressions responsible for them. It should be possible to track source location information for only the most recent potential culprit for each type of coercion failure, combining space-efficient gradual typing with informative error messages.

ACKNOWLEDGMENTS

David Herman is supported by a grant from the Mozilla Corporation. Cormac Flanagan is supported by a Sloan Fellowship and by NSF grant CCR-0341179. Aaron Tomb is also supported by NSF grant CCR-0341179.

REFERENCES

- [1] G. Birtwhistle et al. *Simula Begin*. Chartwell-Bratt Ltd., 1979.
- [2] G. Bracha. Pluggable type systems. In *Workshop on Revival of Dynamic Languages*, October 2004.
- [3] L. Cardelli. Amber. In *Spring School of the LITP on Combinators and Functional Programming Languages*, pages 21–47, 1986.
- [4] L. Cardelli, J. Donahue, L. Glassman, M. Jordan, B. Kalsow, and G. Nelson. Modula-3 report (revised). Technical Report 52, DEC SRC, 1989.
- [5] C. Chambers. *The Cecil Language Specification and Rationale: Version 3.0*. University of Washington, 1998.
- [6] J. Clements. *Portable and high-level access to the stack with Continuation Marks*. PhD thesis, Northeastern University, 2005.

- [7] R. B. de Oliveira. The Boo programming language, 2005.
- [8] Ecma International. *ECMAScript Language Specification*, third edition, December 1999.
- [9] Ecma International. *ECMAScript Edition 4 group wiki*, 2007.
- [10] M. Felleisen and M. Flatt. Programming languages and lambda calculi. Lecture notes online, July 2006.
- [11] R. B. Findler and M. Felleisen. Contracts for higher-order functions. In *International Conference on Functional Programming*, pages 48–59, Oct. 2002.
- [12] C. Flanagan. Hybrid type checking. In *Symposium on Principles of Programming Languages*, pages 245–256, 2006.
- [13] C. Flanagan, S. N. Freund, and A. Tomb. Hybrid types, invariants, and refinements for imperative objects. In *International Workshop on Foundations and Developments of Object-Oriented Languages*, 2006.
- [14] S. E. Ganz, D. P. Friedman, and M. Wand. Trampolined style. In *International Conference on Functional Programming*, pages 18–27, 1999.
- [15] J. J. Garrett. Ajax: A new approach to web applications, 2005.
- [16] J. Gronski, K. Knowles, A. Tomb, S. N. Freund, and C. Flanagan. Sage: Hybrid checking for flexible specifications. In *Scheme and Functional Programming Workshop*, September 2006.
- [17] F. Henglein. Dynamic typing: Syntax and proof theory. *Sci. Comput. Program.*, 22(3):197–230, 1994.
- [18] J. Matthews and R. B. Findler. Operational semantics for multi-language programs. In *POPL '07: Conference record of the 34th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, 2007.
- [19] E. Meijer and P. Drayton. Static typing where possible, dynamic typing when needed. In *Workshop on Revival of Dynamic Languages*, 2005.
- [20] G. Morrisett, M. Felleisen, and R. Harper. Abstract models of memory management. In *International Conference on Functional Programming Languages and Computer Architecture*, pages 66–77, 1995.
- [21] J. G. Siek and W. Taha. Gradual typing for functional languages. In *Scheme and Functional Programming Workshop*, September 2006.
- [22] S. Tobin-Hochstadt and M. Felleisen. Interlanguage migration: From scripts to programs. In *Dynamic Languages Symposium*, October 2006.
- [23] A. K. Wright. *Practical Soft Typing*. PhD thesis, Rice University, Aug. 1998.
- [24] A. K. Wright and M. Felleisen. A syntactic approach to type soundness. *Information and Computation*, 115(1):38–94, 1994.