Into the Core Squeezing Haskell into nine constructors

> Simon Peyton Jones Microsoft Research September 2016

CACM Jan 2016

DOI:10.1145.2844112

Article development led by aCMQUEUE queue.acm.org

We need it, we can afford it, and the time is now.

BY PAT HELLAND

Immutability Changes Everything



GHC Haskell

A very complicated and ill-defined language, with a long user manual, that almost no one understands completely



GHC is big and old

Module	Lines (1992)	Lines (2011)	Increase
Compiler			-
Main	997	11,150	11.2
Parser	1,055	4,098	3.9
Renamer	2,828	4,630	1.6
Type checking	3,352	24,097	7.2
Desugaring	1,381	7,091	5.1
Core tranformations	1,631	9,480	5.8
STG transformations	814	840	1
Data-Parallel Haskell		3,718	1000
Code generation	2913	11,003	3.8
Native code generation		14,138	
LLVM code generation	35.00	2,266	
GHCi		7,474	
Haskell abstract syntax	2,546	3,700	1.5
Core language	1,075	4,798	4.5
STG language	517	693	1.3
C (was Abstract C)	1,416	7,591	5.4
Identifier representations	1,831	3,120	1.7
Type representations	1,628	3,808	2.3
Prelude definitions	3,111	2,692	0.9
Utilities	1,989	7,878	3.96
Profiling	191	367	1.92
Compiler Total	28,275	139,955	4.9
Runtime System			
All C and C code	43,865	48,450	1.10

Figure 1: Lines of code in GHC, past and present

	Module	Lines (1992)	Lines (2011)	Increase
	Compiler			
	Main	997	11,150	11.2
OTTO	Parser	1,055	4,098	3.9
	Renamer	2,828	4,630	1.6
	Type checking	3,352	24,097	7.2
	Desugaring	1,381	7,091	5.1
•	Core tranformations	1.631		Y
10	STG transform t			
			5	
		0STIU		
1 •	Qu			
hio	Qu			
big	Qu			
big	Qu			
big	Qu			2
big	Qu		can	07
big and	Qui La to	ctav	san	e?
big and	how to	stay	san	e?
big and	how to	stay	San 2.692	e?
big and old	how to	stay	SAN 2,692 7,878	e? 2.3 0.9 3.96
big and old	how to	stay 3,111 1,989 191	SAN 2,692 7,878 367	e? 2.3 0.9 3.96 1.92
big and old	how to Profiling Compiler Total	stay stay ^{3,111} 1,989 191 28,275	SAN 2,692 7,878 367 139,955	2.3 0.9 3.96 1.92 4.9
big and old	And the system	stay 5,111 1,989 191 28,275	SAN 2,692 7,878 367 139,955	2.3 0.9 3.96 1.92 4.9



How GHC works

Source language

Typed intermediate language

Massive language Hundreds of pages of user manual Syntax has dozens of data types 100+ constructors

Core 3 types, ypecheck **15 constructors** Desugar Rest of GHC

A typed intermediate language

Haskell	Core (the typed IL)
Big	Small
Implicitly typed	Explicitly typed
Binders typically un-annotated \x. x && y	Every binder is type-annotated \(x:Bool). x && y
Type inference (complex, slow)	Type checking (simple, fast)
Complicated to specify just which programs will type-check	Very simple to specify just which programs are type-correct
Ad-hoc restrictions to make inference feasible	Very expressive indeed; simple, uniform

A typed intermediate language

\x. x && y	\(x:Bool). x && y
Type inference (complex, slow)	Type checking (simple, fast)
Complicated to specify just which programs will type-check	Very simple to specify just which programs are type-correct
Ad-hoc restrictions to make inference feasible	Very expressive indeed; simple, uniform

A typed intermediate language: why?

- 1. Small IL means that analysis, optimisation, and code generation, handle only a small language.
- 2. Type checker ("Lint") for Core is a very powerful internal consistency check on most of the compiler
 - Desugarer must produce well-typed Core
 - Optimisation passes must transform well-typed Core to well-typed Core
- 3. Design of Core is a powerful sanity check on crazy type-system extensions to source language. If you can desugar it into Core, it must be sound; if not, think again.

A typed intermediate language: why?

- 1. Small IL means that analysis, optimication handle only a small GHC is the only production
- 2. Type checker (check on most of
 - Desugarer mus
 - Optimisation pc
- 3. Design of Core is extensions to sou must be sound; if

GHC is the only product. GHC is the only product. compiler that remorselessly pursues this idea of a strongly-typed intermediate language

The design of Core is probably GHC's single most substantial technical achievement

de generation,

ternal consistency

l-typed Core

pe-system 1 into Core, it

Why should you care about Core?

- Because it is jolly interesting
- Because much of it would work for Erlang
- Because it shows theory and practice working together

WHAT SHOULD CORE BE LIKE?

What should Core be like?

- Start with lambda calculus. From "Lambda the Ultimate X" papers we know that lambda is super-powerful.
- But we need a TYPED lambda calculus

Idea:

- start with lambda calculus
- sprinkle type annotations

But:

- Don't want to be buried in type annotations
- Types change as you optimise

Example

 Idea: put type annotations on each binder (lambda, let), but nowhere else

- But: where are 'a' 'b' 'c' bound?
- And: unstable under transformation...

Example

compose :: (b->c) -> (a->b) -> a -> c
compose = λf:b->c. λg:a->b. λx:a.
 let tmp:b = g x
 in f tmp

neg :: Int -> Int isPos :: Int -> Bool

compose isPos neg
= (inline compose:
 f=isPos, g=neg)
Ax:a. let tmp:b = neg x
 in isPos tmp

Now the type annotations are wrong
 Solution: learn from Girard and Reynolds!

System F

compose :: $\forall abc. (b \rightarrow c) \rightarrow (a \rightarrow b) \rightarrow a \rightarrow c$ compose = $\Lambda abc. \lambda f: b \rightarrow c. \lambda g: a \rightarrow b. \lambda x: a.$ let tmp: b = g x in f tmp

Idea: an explicit (big) lambda binds type variables

System F

compose :: $\forall abc. (b \rightarrow c) \rightarrow (a \rightarrow b) \rightarrow a \rightarrow c$ compose = $\Lambda abc. \lambda f:b \rightarrow c. \lambda g:a \rightarrow b. \lambda x:a.$ let tmp:b = g x in f tmp

compose @Int @Int @Bool isPos neg

- = (inline compose: a=Int, b=Int, c=Bool, f=isPos, g=neg) λx:Int. let tmp:Int = neg x in isPos tmp
- Big lambdas are applied to types, just as little lambdas are applied to values
- Now the types stay correct!



Core: GHC's intermediate language

data Expr	
= Var	Var
Lit	Literal
App	Expr Expr
Lam	Var Expr Both term and type lambda
Let	Bind Expr
Case	<pre>Expr Var Type [(AltCon, [Var], Expr)]</pre>
Type	Type Used for type application
data Var = Io	d Name Type Term variable

```
| TyVar Name Kind -- Type/kind variable
```

```
type Kind = Type
data Type = TyVarTy Var
    | LitTy TyLit
    AppTy Type Type
    TyConApp TyCon [Type]
    FunTy Type Type -- Not really necy
    ForAllTy Var Type
```

Core: GHC's intermediate language

```
| CaseExpr Var Type [(AltCon, [Var], Expr)]| TypeType--Used for type application
```

```
type Kind = Type
data Type = TyVarTy Var
    | LitTy    TyLit
    AppTy    Type Type
    | TyConApp TyCon [Type]
    | FunTy    Type Type -- Not really necy
    | ForAllTy Var Type
```

Core: GHC's intermediate language



What's good about System F In our presentation of System F, each variable occurrence is annotated with its type.

Hence every term has a unique type

```
exprType :: Expr -> Type
exprType (Var v) = varType v
exprType (Lam v a) = Arrow (varType v) (exprType a)
...more equations...
```

exprType is pure; needs no "Gamma" argument

Sharing of the Var means that the apparent duplication is not real

What's good about System F?

Type checking (Lint) is fast and easy, because the rules are syntax-directe d



What's good about System F?

Type checking (Lint) is fast and easy, because the rules are syntax-directe d



The syntax of a term encodes its typing derivation

Story so far

- Very small, statically typed language
- Robust to transformations (ie if the term is well typed, then the transformed term is well typed)
- Simple, pure exprType
- Type checking (Lint) is easy and fast

THE CORE PIPELINE



- The Simplifier
 - Inlining
 - Rewrite rules
 - Beta reduction
 - Case of case
 - Case of known constructor
 - etc etc etc
- Specialise overloading
- Float out
- Float in
- Demand, cardinality, and CPR analysis
- Arity analysis
- Call-pattern specialisation (SpecConstr)



```
bash$ ghc -c Foo.hs -ddump-simpl
Result size of Tidy Core
  = {terms: 7, types: 4, coercions: 0}
Foo.f :: GHC.Types.Int -> GHC.Types.Int
[GblId, Arity=1, Str=DmdType]
Foo_f =
 (x apE :: GHC.Types.Int) \rightarrow
   GHC.Num.+ @ GHC.Types.Int
            GHC.Num.$fNumInt
            x apE (GHC.Types.I# 1)
```

```
bash$ ghc -c Foo.hs -dshow-passes
*** Checking old interface for main: Foo:
*** Parser:
*** Renamer/typechecker:
*** Desugar:
Result size of Desugar (after optimization)
  = {terms: 7, types: 4, coercions: 0}
*** Simplifier:
Result size of Simplifier = {terms: 7, types: 4, coercions: 0}
*** Tidy Core:
Result size of Tidy Core = {terms: 7, types: 4, coercions: 0}
*** CorePrep:
Result size of CorePrep = {terms: 9, types: 5, coercions: 0}
*** Stg2Stg:
*** CodeOutput:
*** New CodeGen:
```

```
bash$ ghc -c Foo.hs -ddump-simpl -O
Result size of Tidy Core
   = {terms: 9, types: 5, coercions: 0}
Foo.f :: GHC.Types.Int -> GHC.Types.Int
[GblId, Arity=1, Caf=NoCafRefs,
Str=DmdType <S,1*U(U)>m,
 Unf=Unf{Src=InlineStable, TopLvl=True,
        Arity=1, Value=True,
        ConLike=True, WorkFree=True, Expandable=True,
        Guidance=ALWAYS IF (unsat ok=True,
                          boring ok=False)
        Tmpl= ...]
Foo_f =
  \ (x arV :: GHC.Types.Int) \rightarrow
   case x arV of [Occ=Dead] { GHC.Types.I# x1 aNv ->
   GHC.Types.I# (GHC.Prim.+# x1 aNv 1)
```

bash\$ ghc -c Foo.hs -O -dshow-passes *** Parser: *** Renamer/typechecker: *** Desugar: Result size of Desugar (after optimization) = {terms: 7, types: 4, coercions: 0} *** Simplifier: Result size of Simplifier iteration=1 = {terms: 6, types: 3, coercions: 0} Result size of Simplifier = {terms: 6, types: 3, coercions: 0} *** Specialise: Result size of Specialise = {terms: 6, types: 3, coercions: 0} *** Float out(FOS {Lam = Just 0, Consts = True, PAPs = False}): Result size of Float out(FOS {Lam = Just 0,Consts = True,PAPs = False}) = {terms: 8, types: 4, coercions: 0} *** Float inwards: Result size of Float inwards = {terms: 8, types: 4, coercions: 0} *** Simplifier: Result size of Simplifier iteration=1 = {terms: 12, types: 6, coercions: 0} Result size of Simplifier = {terms: 9, types: 5, coercions: 0} *** Simplifier: Result size of Simplifier = {terms: 9, types: 5, coercions: 0} *** Simplifier: Result size of Simplifier = {terms: 9, types: 5, coercions: 0} *** Demand analysis: Result size of Demand analysis = {terms: 9, types: 5, coercions: 0} *** Worker Wrapper binds: Result size of Worker Wrapper binds = {terms: 9, types: 5, coercions: 0} *** Simplifier: Result size of Simplifier = {terms: 9, types: 5, coercions: 0} *** Float out(FOS {Lam = Just 0, Consts = True, PAPs = True}): Result size of Float out(FOS {Lam = Just 0, Consts = True, PAPs = True}) = {terms: 9, types: 5, coercions: 0} *** Common sub-expression: Result size of Common sub-expression = {terms: 9, types: 5, coercions: 0} *** Float inwards: Result size of Float inwards = {terms: 9, types: 5, coercions: 0} *** Simplifier: Result size of Simplifier = {terms: 9, types: 5, coercions: 0} *** Tidy Core: Result size of Tidy Core = {terms: 9, types: 5, coercions: 0} *** CorePrep: Result size of CorePrep = {terms: 12, types: 6, coercions: 0} *** Sta2Sta: *** CodeOutput: *** New CodeGen:

CORE OPTIMISATIONS

A hail of bullets



- The Simplifier
 - Inlining
 - Rewrite rules
 - Beta reduction
 - Case of case
 - Case of known constructor
 - etc etc etc
- Specialise overloading
- Float out
- Float in
- Demand, cardinality, and CPR analysis
- Arity analysis
- Call-pattern specialisation (SpecConstr)

The mighty simplifier

- Simplifier performs a raft of local optimisations
 - Inlining
 - Beta reduction
 - Rewrite rules
 - Case of known constructor
 - Case of case
 - etc etc etc
- These transformations cascade

LET: Inlining

Inlining is the key transformation

Replace x by x's definition

let x = <rhs> in ...x...x...
===>
let x = <rhs> in ...x..<rhs>...

Lots of heuristics here

Issues:

- May duplicate code <rhs> How often does 'x' occur? How big is <rhs>?
- May duplicate work let x = factorial y in ... Is <rhs> a value or not?
- Which occurrences to substitute for. All? Some? None?
LET: Inlining

Works even if x's definition is in another module, or another package

module Prelude where
null :: [a] -> Bool
null [] = True
null (x:xs) = False
not :: Bool -> Bool
not True = False
not False = True

module X where import Prelude notNull xs = not (null xs)

Prelude.hi Gives type signature and unfolding for null, not

X.hi

Beta reduction is easy: make a 'let'

LAMBDA: Beta reduction

(\x. e1) e2 ===> let x = e2 in e1

reduction Note that: beta reduction produces a 'let'; now inlining decisions take over again.

> Other compilers: inline function call GHC: inline occurrence of function; then beta reduce

APPLY:

Lots of optimisations of the form fun arg₁ ... arg_n ==> result

Rewrite A large number are built-in to GHC; rules very easy to extend this list +#



Simplifier: applies the rules exhaustively

 Simple "phase system" to control order of application. Apply phase-2 rules; then phase-1 rules; then phase-0 rules.

Users can write more RULES

Rewrite rules

- Works across modules
- Effectively: extends the compiler with library-specific optimisation rules!
- Compiler does not (cannot) check correctness; that's up to the programmer.

Compiler writes some new rewrite rules itself

Compiler-g enerated rewrite rules



Compiler-g enerated rewrite rules Again, two stages

Definition site: generate the specialised function, and the rewrite rule

Call sites: use the rewrite rule, call site by call site, to optimise the program

Works beautifully across modules

Case of known constructor

CASE: case rules

case True of { False -> e1; True -> e2 }
 ===>
e2

Savvy to enclosing let-binding

let x = a:as in ...(case x of
 [] -> e1
 (p:ps) -> e2)...
===>
let x = a:xs in ...(let p=a; ps=as in e2)...

Case of known constructor

CASE: case rules

case True of { False -> e1; True -> e2 } ===> e2

...and to enclosing case

```
case x of
  True -> ...(case x of {True -> e1; False -> e2 })...
  False -> e3
===>
case x of
  True -> ...e1...
False -> e3
```

	<pre>notNull xs = not (null xs)</pre>	
ASE: ase of	= case (null xs) of True -> False False -> True	I
ase	<pre>= case (case xs of [] -> True (p:ps) -> False) of True -> False</pre>	Ţ
	False -> True	

```
module Prelude where
null :: [a] -> Bool
null [] = True
null (x:xs) = False
not :: Bool -> Bool
not True = False
not False = True
```



A worry

=	case (c	ase xs of	
		[] -> True	
		(p:ps) -> False)	of
	True	-> BIG1	
	False	e -> BIG2	

= case xs of
[] -> case True of
 True -> BIG1
 False -> BIG2
(p:ps) -> case False of
 True -> BIG1
 False -> BIG1
 False -> BIG2

Duplicates arbitrary amounts of code :-(

Join points (continuations)



No duplication :-)



Join points



Arbitrary patterns

=	case	(case	xs	of		
		[]		->	True	
		(p	:ps)) ->	False)	of
	Jus	st x	-> I	BIG1		
	Not	hing	-> I	BIG2		

Simply abstract over the pattern bound variables

= let j1 x = BIG1
 j2 () = BIG2
in case xs of
[] -> case True of
 Just x -> j1 x
 False -> j2 ()
(p:ps) -> case False of
 True -> j1 x
 False -> j2 ()

 Works fine for existentials, GADTs Join points are like control-flo w labels

printf("yes") ;
if blah then { stuff }
 else { other stuff } ;
printf("yes");

Let bindings allocate a thunk or data constructor

Join points allocate nothing!

"Calling" a join
 point = adjust
 stack pointer
 and goto

Hot news: Sequent Core



Bad bad bad!

- 'j' is no longer a join point
- 2. The R1/R2 case does not scrutinise E1



Sequent Core Based on Sequent Calculus

But in "direct style"

 Turns out to be very modest variant of Core; easy to implement

Some amazing new fusion happens

THINGS I HAVE NOT TOLD YOU

More in Core

Module-level transformations

- Strictness analysis
- Specialisation
- Float in, float out, lambda lifting
- Static argument transformation
- All of Haskell's exotic type system is in Core
 - Existentials
 - GADTs
 - Newtypes, coercions, casts
 - Kind polymorphism
 - Type level functions
 - Roles

Still only nine constructors

data Expr	
= Var	Var
Lit	Literal
App	Expr Expr
Lam	Var Expr
Let	Bind Expr
Case	<pre>Expr Var Type [(AltCon, [Var], Expr)]</pre>
Type	Туре
Coercion	Coercion
Cast Exp	pr Coercion
data Var = Io	d Name Type
T ₃	yVar Name Kind
Co	oVar Name Type Type
data Bind = I	Rec [(Var,Expr)] NonRec Var Expr

Conclusion

- Typed lambda calculus is totally amazing
 Little here is Haskell specific
- Cross-module inlining (SO IMPORTANT) might be problematic for Erlang
- 25 years on, still in active development:
 - Sequent Core
 - Strict Core

EXISTENTIALS, GADTS, AND COERCIONS

Existentials



T1 :: ∀ab. b -> (b -> a) -> T a ≡ ∀a. (∃b.(b, b->a)) -> T a

Reminder data types and pattern matching



Just 4 :: Maybe Int Just True :: Maybe Bool Nothing :: Maybe Int Nothing :: Maybe Bool Just (Just True) :: Maybe (Maybe Bool)

A typical evaluator

```
data Term = Lit Int
| Succ Term
| IsZero Term
| If Term Term Term
```

```
data Value = VInt Int | VBool Bool
```

```
eval :: Term -> Value
eval (Lit i) = VInt i
eval (Succ t = case eval t of { VInt i -> VInt (i+1) }
eval (IsZero t) = case eval t of { VInt i -> VBool (i==0) }
eval (If b t1 t2)= case eval b of
VBool True -> eval t1
VBool False -> eval t2
```

Richer data types

What if you could define data types with richer return types? Instead of this:

data Term where Lit :: Int -> Term Succ :: Term -> Term IsZero :: Term -> Term If :: Term -> Term -> Term -> Term

we want this:

data Term a where Lit :: Int -> **Term Int** Succ :: Term Int -> **Term Int** IsZero :: Term Int -> **Term Bool** If :: Term Bool -> Term a -> Term a -> Term a

Richer data types

data Term a where Lit :: Int -> Term Int Succ :: Term Int -> Term Int IsZero :: Term Int -> Term Bool If :: Term Bool -> Term a -> Term a -> Term a

Major payoff: rules out ill-typed Terms

Lit 4 :: Term Int IsZero (Lit 4) :: Term Bool If (Lit 4) :: *** TYPE ERROR *** IF (IsZero (Lit 4)) :: Term a -> Term a -> Term a



Now you can write a cool typed evaluator



Evaluator is easier to read and write

Evaluator is more efficient too

GADTs

Generalised Algebraic Data Types (GADTs):

Single idea: allow arbitrary return type for constructors, provided outermost type constructor is still the type being defined

data Term a where Lit :: Int -> Term Int Succ :: Term Int -> Term Int IsZero :: Term Int -> Term Bool If :: Term Bool -> Term a -> Term a -> Term a

GADTs in Core?



Solution to both problems: EVIDENCE



Evidence

T1 :: \forall a. (a~Bool) -> Bool -> T a

- Any application of T1 must supply evidence T1 σ e1 e2 where e1 : (σ~Bool), e2 : Bool
- Here e1 is a value that denotes evidence that o=Bool
- And any pattern match on T1 gives access to evidence
 case s of { T1 (c:σ~Bool) (y:Bool) -> ... }
 where s : Tσ



Modifications to Core

data Expr	
= Var	Var
Lit	Literal
App	Expr Expr
Lam	Var Expr
Let	Bind Expr
Case	<pre>Expr Var Type [(AltCon, [Var], Expr)]</pre>
Type	Туре
Coercion	Coercion Used for coercion apps
Cast Exp	pr Coercion Type-safe cast
data Var = Io	d Name Type Term variable vVar Name Kind Type variable
Co	oVar Name Type Type Coercion var

Evidence terms

T1 :: \forall a. (a~Bool) -> Bool -> T a

- Consider the call: T1 Bool <Bool> True : T Bool
- Here <Bool> : Bool ~ Bool
 γ ::= <T> | ...
- Can I call T1 Char γ True : T Char?
- No: that would need (y : Char ~ Bool) and there are no such terms y
Composing evidence terms



 $\mathbf{Y} ::= \langle \mathbf{T} \rangle$ | Sym \mathbf{Y} | ... If $\mathbf{Y} : \mathbf{T} \sim \sigma$ then Sym $\mathbf{Y} : \sigma \sim \mathbf{T}$

Composing evidence terms



$$\mathbf{Y} ::= \langle \mathbf{T} \rangle | \operatorname{Sym} \mathbf{Y} | \mathbf{T} \mathbf{Y}_1 \dots \mathbf{Y}_n | \dots$$

$$If \mathbf{Y}_i : \mathbf{T}_i \sim \sigma_i$$

$$then \mathbf{T} \mathbf{Y}_1 \dots \mathbf{Y}_n : \mathbf{T} \mathbf{T}_1 \dots \mathbf{T}_n \sim \mathbf{T} \sigma_1 \dots \sigma_n$$

Evidence terms



Cost model

- Coercions are computationally irrelevant
- Coercion abstractions, applications, and casts are erased at runtime

Bottom line

- Just like type abstraction/application, evidence abstraction/application provides a simple, elegant, consistent way to
 - express programs that use local type equalities
 in a way that is fully robust to program transformation
 and can be typechecked in an absolutely straightforward way
- Cost model: coercion abstractions, applications, and casts are erased at runtime



newtypes

Haskell

newtype Age = MkAge Int

bumpAge :: Age -> Int -> Age bumpAge (MkAge a) n = MkAge (a+n)

- No danger of confusing Age with Int
- Type abstraction by limiting visibility of MkAge
- Cost model: Age and Int are represented the same way

In Core?

newtype Age = MkAge Int

bumpAge :: Age -> Int -> Age bumpAge (MkAge a) n = MkAge (a+n)

- Newtype constructor/pattern matching turn into casts
- (New) Top-level axiom for equivalence between Age and Int
- Everything else as before



Axioms can be parameterised, of course
No problem with having a polytype in s~t





WRAP UP

Wrap up

- Many more aspects not covered in this talk
 - Roles: nominal and representational equality
 - Optimising coercions
 - "Closed" type families with non-linear patterns, and proving consistency thereof

```
type family Eq a b where
Eq a a = True
Eq a b = False
```

- Heterogeneous equalities; coercions at the type level
- Collapsing types and kinds: * :: *

Wrap up

- Small, statically typed intermediate language
 - Sanity
 - Powerful, general optimisations
 - Lint
- Main "new" idea: programs manipulate evidence along with types and values. This single idea in Core explains multiple source-language concepts:
 GADTs
 - Newtypes
 - Type and data families (both open and closed)

Great example of theory into practice