

# Cardinality Analysis and its Applications: from Glasgow Haskell Compiler to Sharded Blockchains

Ilya Sergey

YaleNUSCollege



# Part I

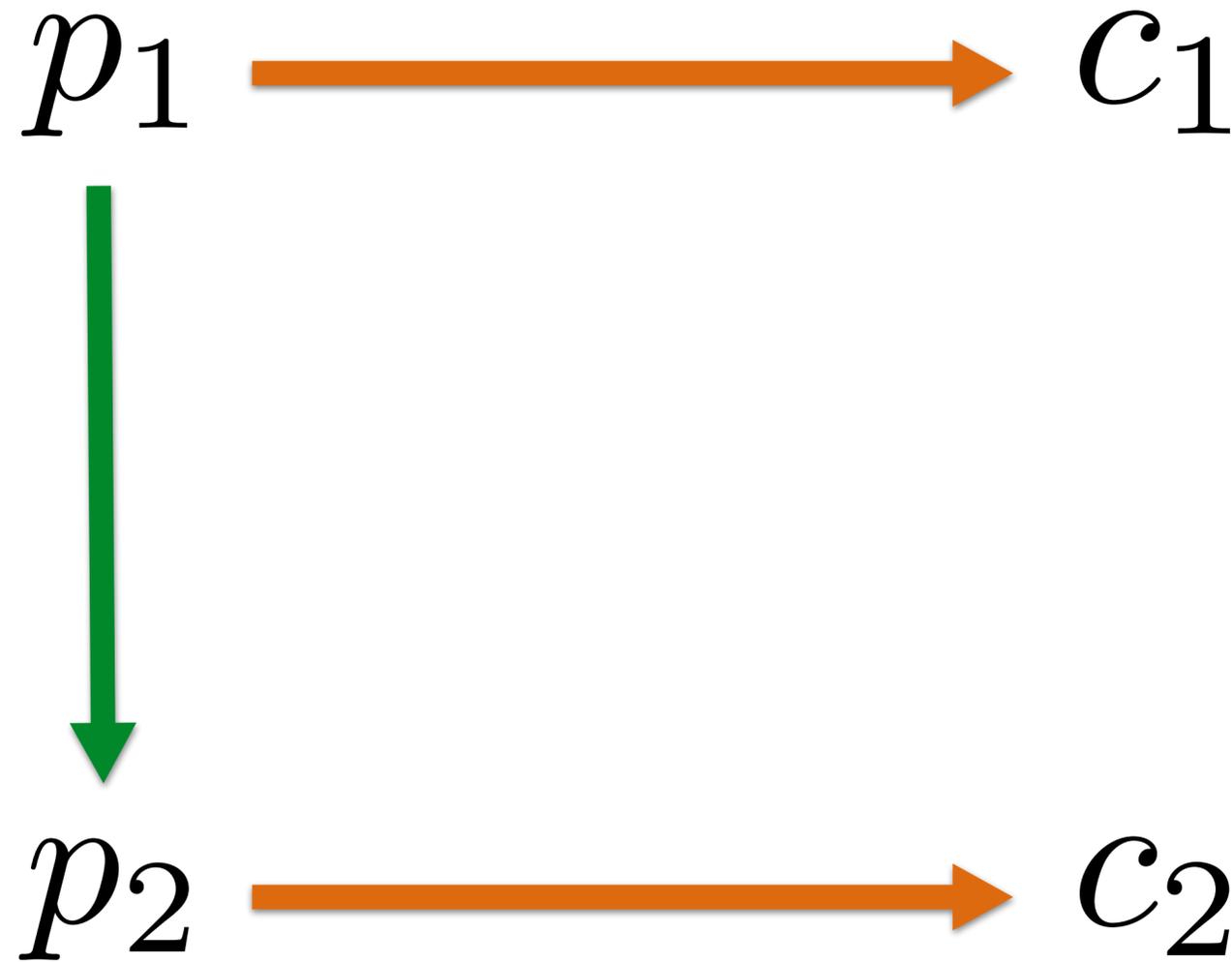
# Cardinality Analysis for Haskell Programs

Joint work with Dimitrios Vytiniotis and Simon Peyton Jones, presented at POPL'14

# Optimising compiler in a nutshell

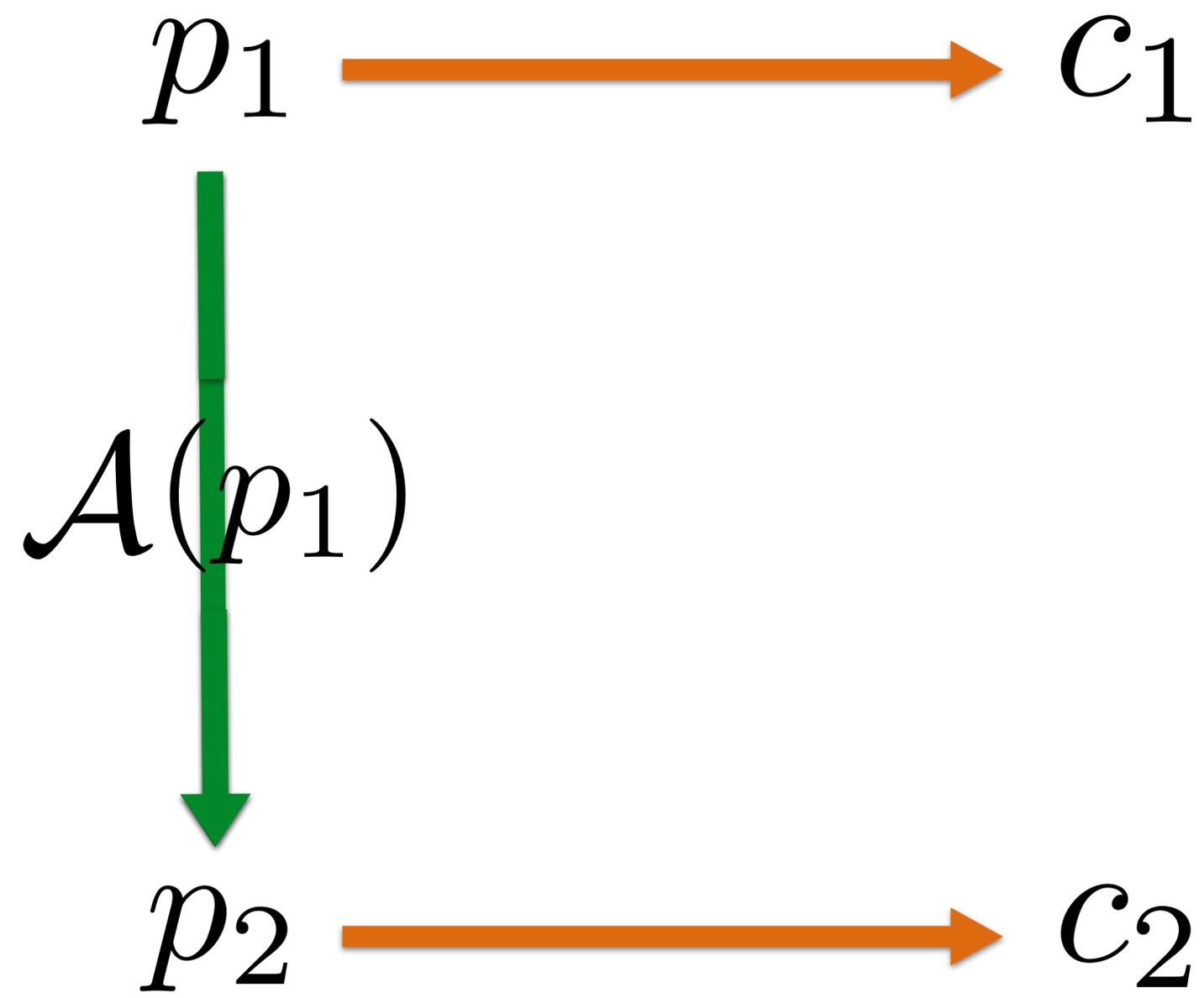
*p*<sub>1</sub>

# Optimising compiler in a nutshell



$$C_2 \cong C_1$$

$$\mathcal{C}_2 \preceq \mathcal{C}_1$$



# Annotating static analysis

*A*

*A story of*  
three program optimisations

# Optimisation 1

f1, f2 :: [Int] -> Int

*Which function is  
better to run?*

Better

f1 xs = let ys = map costly xs  
in **squash** (\n -> sum (map (+ n) ys))

*if invoked more than once by squash*

f2 xs = **squash** (\n -> sum (map (+ n) (map costly xs)))

Better

*if invoked at most once by squash*

# Optimisation 1

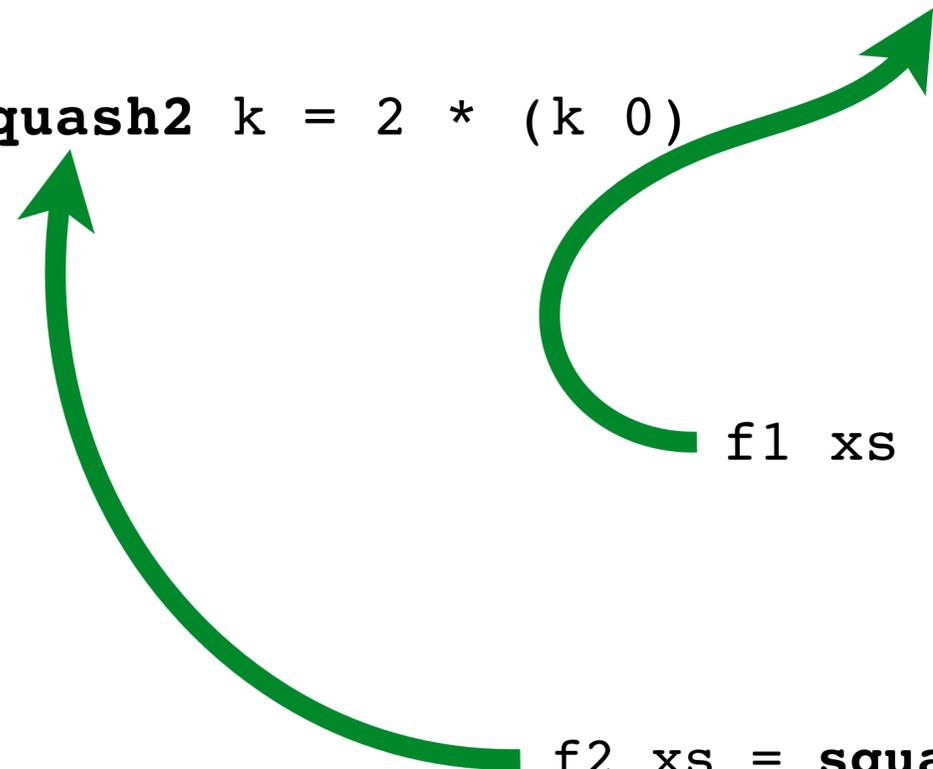
```
squash1, squash2 :: (Int -> Int) -> Int
```

```
squash1 k = sum (map k [1..10])
```

```
squash2 k = 2 * (k 0)
```

```
f1 xs = let ys = map costly xs  
         in squash (\n. sum (map (+ n) ys))
```

```
f2 xs = squash (\n. sum (map (+ n) (map costly xs)))
```



Need to know:  
*how many times*  
a function is called.

(call cardinality)

# Optimisation 2

*“worker-wrapper” split*

```
f x = case x of (p,q) -> <tbody>
```

# Optimisation 2

*“worker-wrapper” split*

*“wrapper”, usually inlined on-site*

`f x = case x of (p,q) -> fw p q`

`fw p q = <cbody>`

*“worker”*

# Optimisation 2

*“worker-wrapper” split*

What if *q* is *never* used in `<cbody>`?

`f x = case x of (p,q) -> fw p`

`fw p = <cbody>`

Don't have to pass *q* to *fw*!

Which parts of  
a data structure are  
certainly *not* used?

(absence)

# Optimisation 3

*smart memoisation*

```
f :: Int -> Int -> Int
f x c = if x > 0 then c + 1 else
        if x == 0 then 0      else c - 1

g y = f y (costly y)
```

Will be used exactly once:  
no need to memoize!

Which parts  
of a data structure  
are used *no more than once*?

(think cardinality)

# Cardinality Analysis

- Call cardinality
- Absence
- Thunk cardinality

# Usage demands

*(how a value is used)*

call demand

Usage demands

$$d ::= \boxed{C^n(d)} \mid U(d_1^\dagger, d_2^\dagger) \mid U$$

Cardinality demands

$$d^\dagger ::= A \mid n * d$$

Usage cardinalities

$$n ::= 1 \mid \omega$$

tuple demand

Usage demands

$$d ::= C^n(d) \mid U(d_1^\dagger, d_2^\dagger) \mid U$$

Cardinality demands

$$d^\dagger ::= A \mid n * d$$

Usage cardinalities

$$n ::= 1 \mid \omega$$

general demand

Usage demands  $d ::= C^n(d) \mid U(d_1^\dagger, d_2^\dagger) \mid \boxed{U}$

Cardinality demands  $d^\dagger ::= A \mid n * d$

Usage cardinalities  $n ::= 1 \mid \omega$

Usage demands

$$d ::= C^n(d) \mid U(d_1^\dagger, d_2^\dagger) \mid U$$

absent value

Cardinality demands

$$d^\dagger ::= \boxed{A} \mid n * d$$

Usage cardinalities

$$n ::= 1 \mid \omega$$

Usage demands

$$d ::= C^n(d) \mid U(d_1^\dagger, d_2^\dagger) \mid U$$

used at most  $n$  times

Cardinality demands

$$d^\dagger ::= A \mid \boxed{n * d}$$

Usage cardinalities

$$n ::= 1 \mid \omega$$

# Usage Types

*(how a function uses its arguments)*

wurble1 ::  $\omega * U \rightarrow C^\omega(C^1(U)) \rightarrow \bullet$

wurble1 a g = g 2 a + g 3 a

$$\text{wurb1e1} \quad :: \quad \omega * U \rightarrow \boxed{C^\omega(C^1(U))} \rightarrow \bullet$$

$$\text{wurb1e1} \quad a \quad g = \boxed{g}^2 \quad a + \boxed{g}^3 \quad a$$

wurble2 ::  $\omega * U \rightarrow C^1(C^\omega(U)) \rightarrow \bullet$

wurble2 a g = sum (map (g a) [1..1000])

wurble2 ::  $\omega * U \rightarrow C^1(C^\omega(U)) \rightarrow \bullet$

wurble2 a g = sum (map (g) a) [1..1000])

$f :: 1 * U(1 * U, A) \rightarrow \bullet$

$f\ x = \text{case } x \text{ of } (p, q) \rightarrow p + 1$

Usage type  
depends on a usage context!

*(result demand determines argument demands)*

# Two Types of Modular Program Analyses

- Forward analysis
  - “Run” the program with *abstract* input and infer the *abstract* result;
  - Examples: sign analysis, interval analysis, type checking/inference.
- Backwards analysis
  - From the expected *abstract* result of the program infer the abstract values of its inputs.

# Backwards Analysis

Infers demand type basing on a context

$$P \vdash \blacktriangleright e \downarrow d \Rightarrow \langle \tau ; \varphi \rangle$$

$$P \mapsto e \downarrow d \Rightarrow \langle \tau ; \varphi \rangle$$

- $P$  - *signature environment*, maps some of free variables of  $e$  to their demand signatures (*i.e.*, keeps some contextual information)
- $d$  - *usage demand*, describes the degree to which  $e$  is evaluated
- $\tau$  - *demand type*, usages that  $e$  places on its arguments
- $\varphi$  - *fv-usage*, usages that  $e$  places on its free variables

$C^1(U)$ 

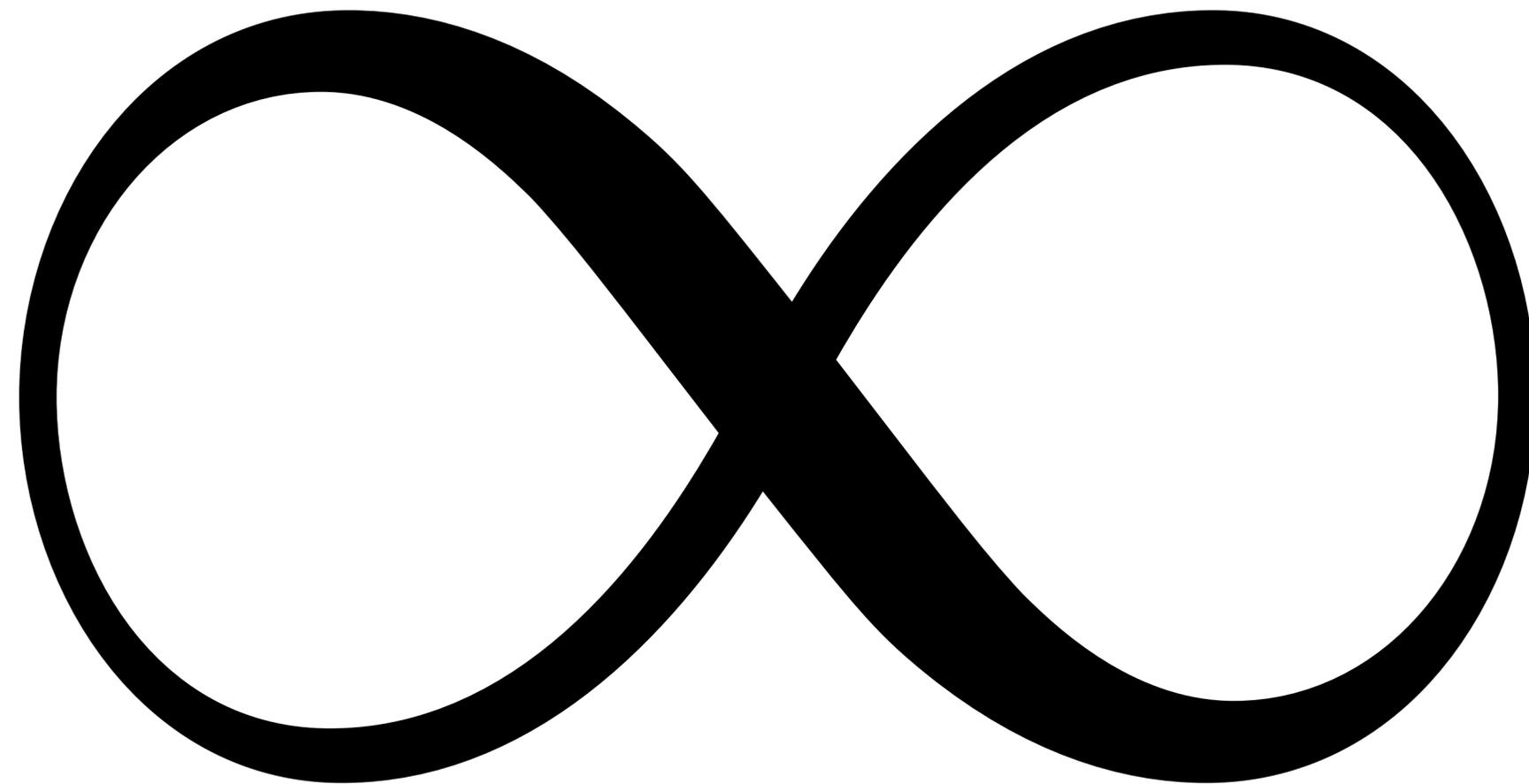
$e = \lambda x . \mathbf{case} \ x \ \mathbf{of} \ (p, q) \rightarrow (p, f \ \mathbf{True})$

$$e = \lambda x . \text{case } x \text{ of } (p, q) \rightarrow (p, \boxed{f \text{ True}})$$

$$\epsilon \vdash e \downarrow C^1(U) \Rightarrow \underbrace{\langle 1 * U(\omega * U, A) \rightarrow \bullet \rangle}_{\mathcal{T}} ; \underbrace{\{\boxed{f \mapsto 1 * C^1(U)}\}}_{\varphi}$$

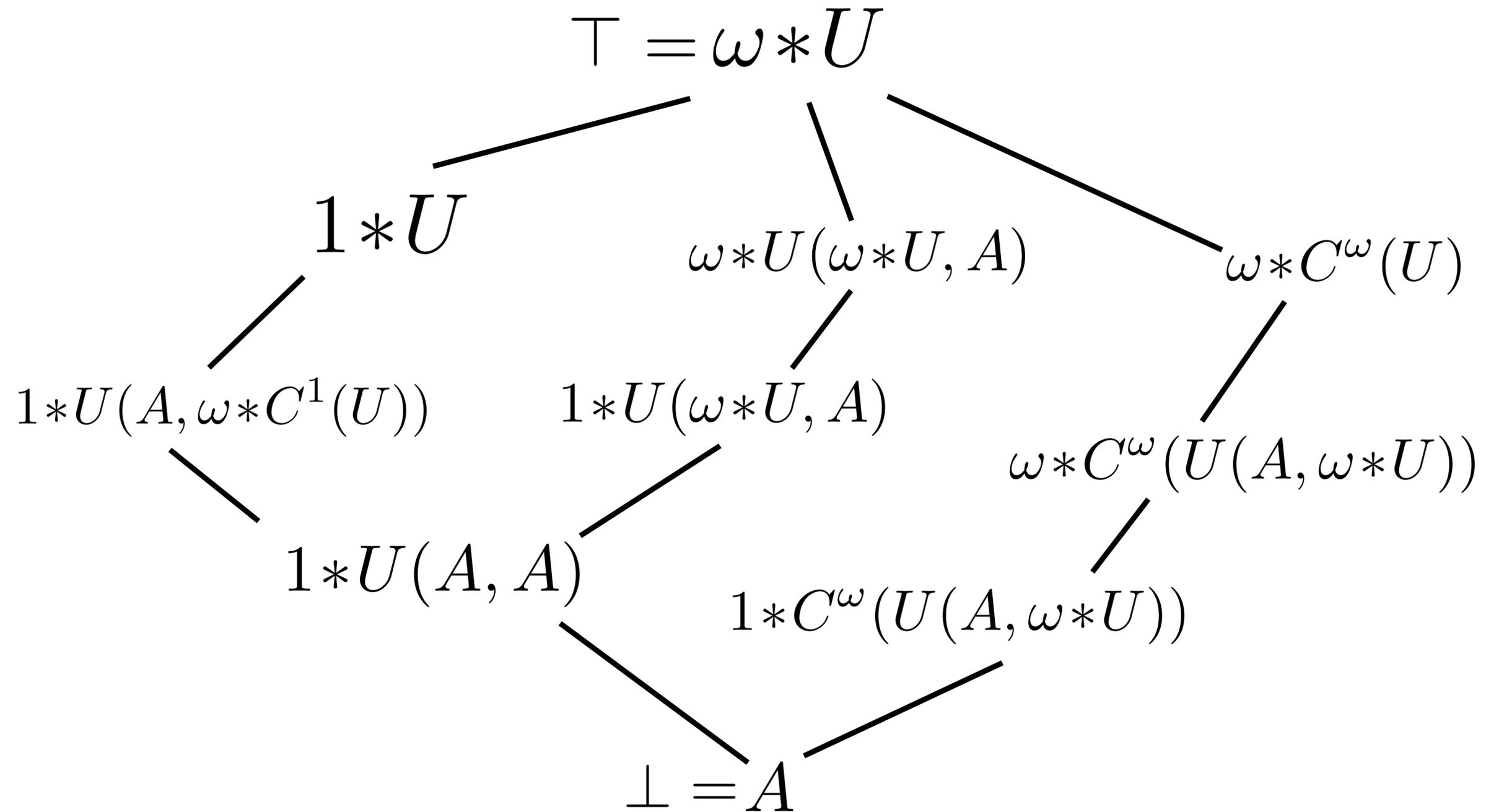
Each function is a  
backwards demand transformer  
it transforms a *context* demand to  
*argument* demands and *fv*-demands.

How many context demands are there?



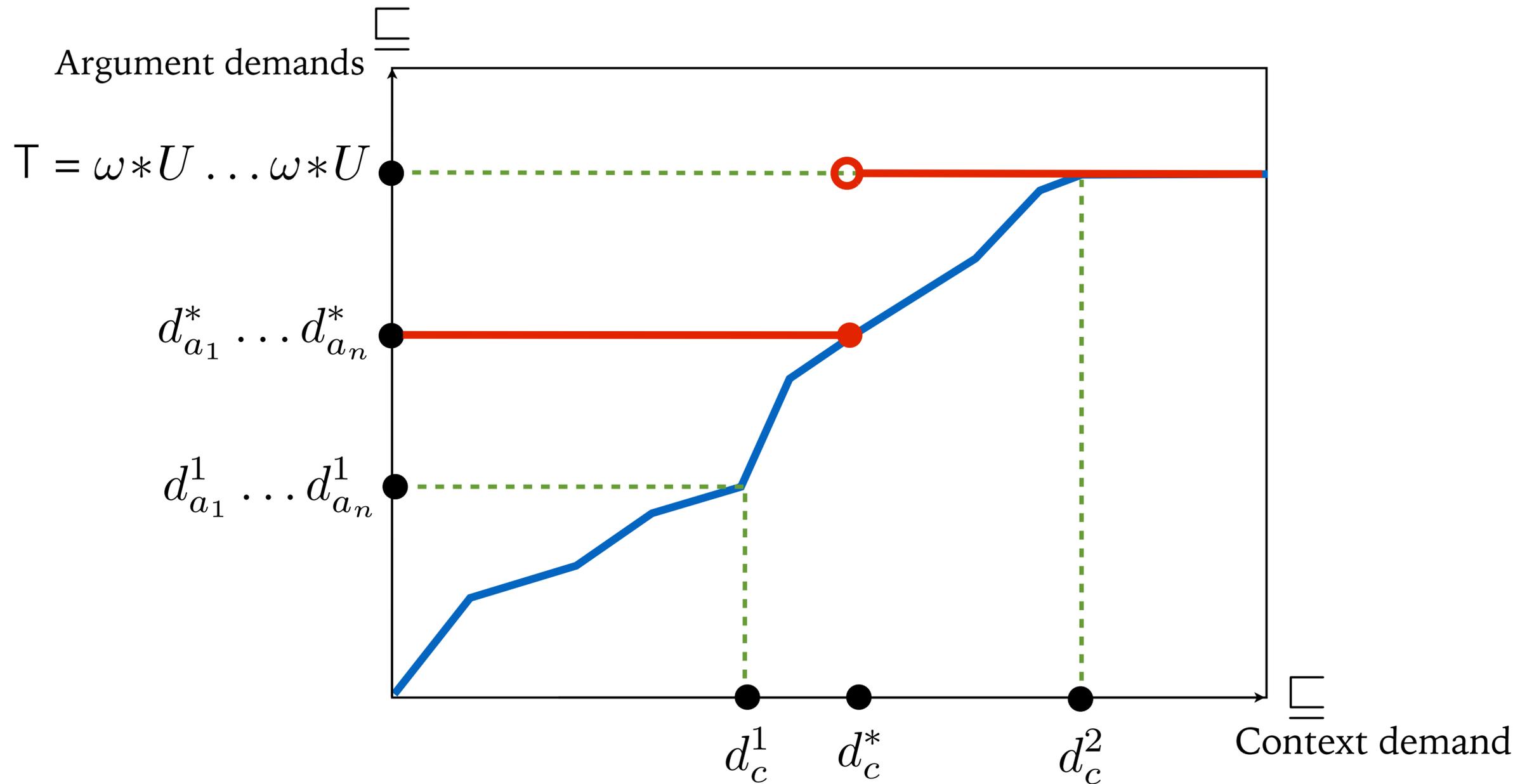
We cannot compute best argument demands  
for *all* contexts:  
need to *approximate*.

# Demand Lattice



Each function is  
a **monotone** backwards demand transformer.

# Exploiting demand monotonicity



# Analysis-based annotations

$$P \mapsto e \downarrow d \Rightarrow \langle \tau ; \varphi \rangle$$

# Elaboration

$$P \vdash e \downarrow d \Rightarrow \langle \tau ; \varphi \rangle \rightsquigarrow e$$

- **let**-bindings in **e** are annotated with  $m \in \{0, 1, \omega\}$  to indicate how often the let-binding is evaluated;
- Each Lambda  $\lambda^n x . e_1$  in **e** carries an annotation  $n \in \{1, \omega\}$  to indicate how often the lambda is called.

$\epsilon \vdash \text{let } f = \lambda x. \lambda y. x \text{ True in } f \ p \ q \ \downarrow C^1(U)$

$\Rightarrow \langle \bullet; \{p \mapsto 1 * C^1(U), q \mapsto A\} \rangle$

$\rightsquigarrow$

$\text{let } f \stackrel{1}{=} \lambda^1 x. \lambda^1 y. x \text{ True in } f \ p \ q$

Soundness

# Restricted operational semantics

(makes sure that the annotations are respected)

Annotating  
cardinality  
analysis

*produces well-typed  
programs*

Type and effect  
system

*annotated programs  
do not get stuck*

*progress and preservation*

Restricted  
operational  
semantics

# Small-Step CBN Machine

Sestoft:JFP97

$$e_1$$

# Small-Step CBN Machine

$$\langle H_1, e_1, S_1 \rangle$$

# Small-Step CBN Machine

$$\langle H_1, e_1, S_1 \rangle \longrightarrow \dots \longrightarrow \langle H_n, e_n, S_n \rangle$$

$$e_1 \rightsquigarrow e_1$$

# Erasing Annotations

$$e_1^{\#} = e_1$$

$\langle H_1, e_1, S_1 \rangle$

# Restricted CBN Machine

$$\langle H_1, e_1, S_1 \rangle \hookrightarrow \dots \hookrightarrow \langle H_n, e_n, S_n \rangle$$

- *1*-annotated lambdas can be called at most once;
- *1*-annotated bindings can be used only once;
- *0*-annotated bindings cannot be used at all.

# Soundness Theorem

An analysis-annotated program  
behaves *the same way* under restricted semantics  
as the original program  
under the normal semantics.

# Soundness Theorem

If  $\epsilon \vdash \triangleright e_1 \downarrow U \Rightarrow \langle \tau, \epsilon \rangle \rightsquigarrow e_1$   
and  $\langle \epsilon ; e_1 ; \epsilon \rangle \xrightarrow{k} \langle H ; e_2 ; S \rangle$

then

$$\langle \epsilon ; e_1 ; \epsilon \rangle \hookrightarrow_{\exists}^k \langle \mathbf{H} ; e_2 ; \mathbf{S} \rangle$$

such that  $\langle \mathbf{H}^{\sharp}, e_2^{\sharp}, \mathbf{S}^{\sharp} \rangle = \langle H, e_2, S \rangle$

# Cardinality-Enabled Optimisations

# 1. Let-in floating optimisation

$\text{let } z \stackrel{m_1}{=} e_1 \text{ in } (\text{let } f \stackrel{m_2}{=} \lambda^1 x . e \text{ in } e_2)$

$\text{let } z \stackrel{m_1}{=} e_1 \text{ in } (\text{let } f \stackrel{m_2}{=} \lambda^1 x . e \text{ in } e_2)$

$\implies \text{let } f \stackrel{m_2}{=} \lambda^1 x . (\text{let } z \stackrel{m_1}{=} e_1 \text{ in } e) \text{ in } e_2;$

for any  $m_1, m_2$  and  $z \notin FV(e_2)$ .

# Improvement Theorem 1

Let-in floating  
*does not* increase the number  
of execution steps.

# Improvement Theorem 1

For any  $H$  and  $S$ , if

$$\langle H ; \text{let } z \stackrel{m}{=} e_1 \text{ in } (\text{let } f \stackrel{m_1}{=} \lambda^1 x . e \text{ in } e_2) ; S \rangle \Downarrow^{\circlearrowleft N}$$

and  $z \notin FV(e_2)$

then

$$\langle H ; \text{let } f \stackrel{m_1}{=} \lambda^1 x . (\text{let } z \stackrel{m}{=} e_1 \text{ in } e) \text{ in } e_2 ; S \rangle \Downarrow^{\leq N}$$

where  $\sigma \Downarrow^N$  means “terminates in  $N$  steps”.

## 2. Smart Execution

$e_1$

# Optimised CBN Machine

$$\langle H_1, e_1, S_1 \rangle \Longrightarrow \dots \Longrightarrow \langle H_n, e_n, S_n \rangle$$

- *1*-annotated bindings are *not* memoised;
- *0*-annotated bindings are *skipped*.

# Improvement Theorem 2

Optimising semantics  
works *faster* on elaborated expressions  
and produces coherent results.

# Improvement Theorem 2

If  $\epsilon \mapsto e_1 \downarrow U \Rightarrow \langle \tau, \epsilon \rangle \rightsquigarrow e_1$

and  $\langle \epsilon ; e_1 ; \epsilon \rangle \xrightarrow{k} \langle H ; e_2 ; S \rangle$

then

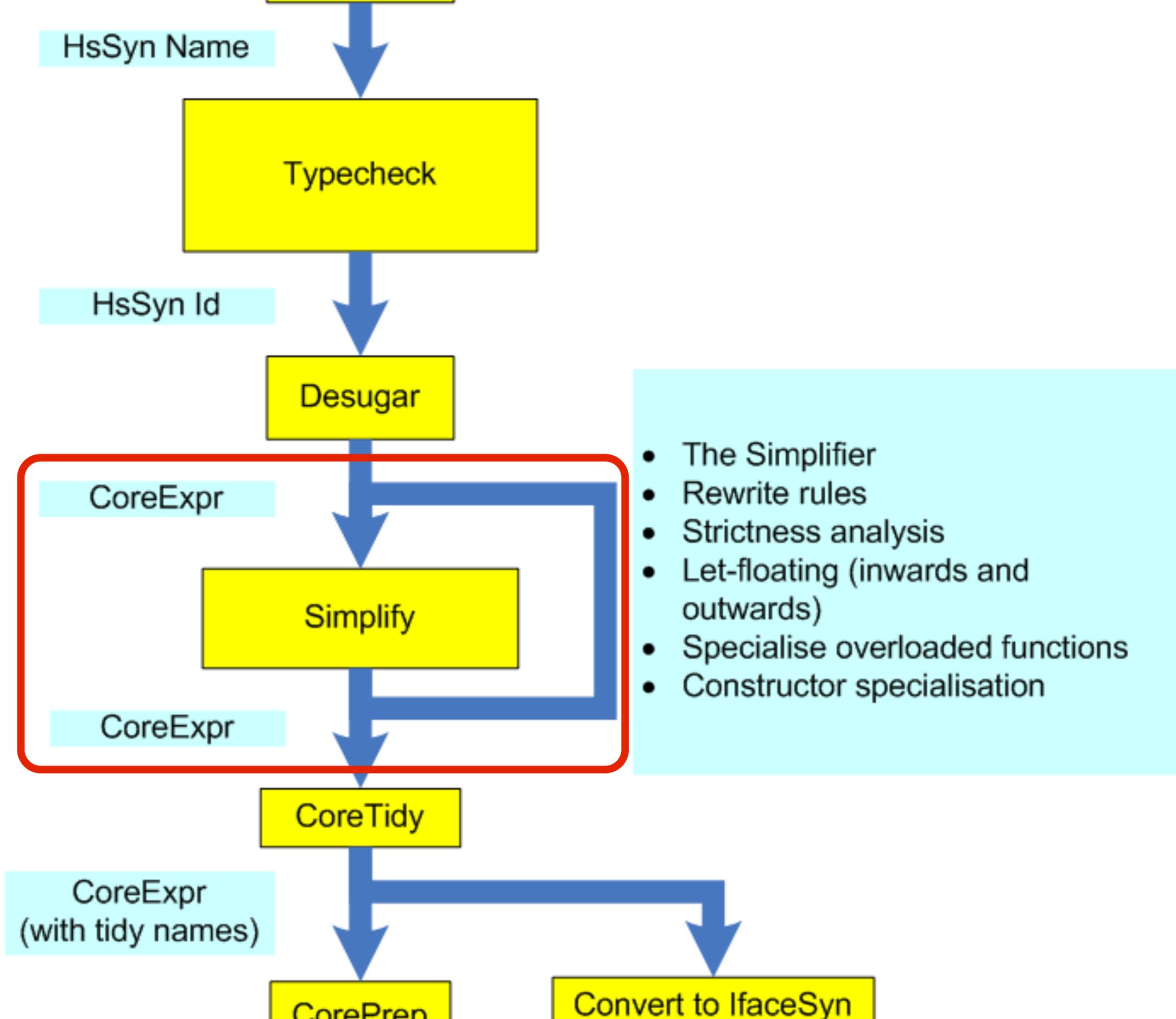
$\langle \epsilon, e_1, \epsilon \rangle \Longrightarrow^m \langle gc(H), e_2, gc(S) \rangle$

such that  $m \leq k$

and  $\langle H^\sharp, e_2^\sharp, S^\sharp \rangle = \langle H, e_2, S \rangle$

Implementation  
and  
Evaluation





# GHC Core

- A tiny language, to which Haskell sources are de-sugared;
- Based on explicitly typed System F with type equality coercions;
- Used as a base platform for analyses and optimisations;
- All names are fully-qualified;
- **if-then-else** is compiled to **case-expressions**;
- Variables have additional metadata;
- Type class constraints are compiled into record parameters.

# Core Syntax

```
data Expr b
  = Var      Id
  | Lit      Literal
  | App      (Expr b) (Expr b)
  | Lam      b (Expr b)
  | Let      (Bind b) (Expr b)
  | Case     (Expr b) b Type [Alt b]
  | Cast     (Expr b) Coercion
  | Tick     (Tickish Id) (Expr b)
  | Type     Type
  | Coercion Coercion

data Bind b = NonRec b (Expr b)
             | Rec [(b, (Expr b))]

type Alt b = (AltCon, [b], Expr b)

data AltCon
  = DataAlt DataCon
  | LitAlt  Literal
  | DEFAULT
```

# How to See Core

## Desugared Core

```
> ghc -ddump-ds Program.hs
```

## Core with Strictness Annotations

```
> ghc -O2 -ddump-stranal Program.hs
```

## Core after Worker/Wrapper Split

```
> ghc -O2 -ddump-worker-wrapper Program.hs
```

# Try it on

```
module Program where
```

```
squash f = f 42
```

```
costly x = product [1..x]
```

```
foo xs = squash (\n -> sum (map (+ n) (map costly xs)))
```

- The analysis and optimisations are implemented in Glasgow Haskell Compiler (GHC v7.8 and newer):  
<http://github.com/ghc/ghc>
- Added 250 LOC to 140 KLOC compiler;
- Runs simultaneously with the strictness analyser;
- Evaluated on
  - **nofib** benchmark suite,
  - various **hackage** libraries,
  - the Benchmark Game programs,
  - GHC itself.

# Results on nofib

Program	Synt. $\lambda^1$	Synt. Thnk <sup>1</sup>	RT Thnk <sup>1</sup>
anna	4.0%	7.2%	2.9%
bspt	5.0%	15.4%	1.5%
cacheprof	7.6%	11.9%	5.1%
calendar	5.7%	0.0%	0.2%
constraints	2.0%	3.2%	4.5%
... and 72 more programs			
Arithmetic mean	10.3%	12.6%	5.5%

\* as linked and run with libraries

# Results on nofib

Program	Allocs		Runtime	
	No hack	Hack	No hack	Hack
anna	-2.1%	-0.2%	+0.1%	-0.0%
bspt	-2.2%	-0.0%	-0.0%	+0.0%
cacheprof	-7.9%	-0.6%	-6.1%	-5.0%
calendar	-9.2%	+0.2%	-0.0%	-0.0%
constraints	-0.9%	-0.0%	-1.2%	-0.2%
... and 72 more programs				
Min	-95.5%	-10.9%	-28.2%	-12.1%
Max	+3.5%	+0.5%	+1.8%	+2.8%
Geometric mean	-6.0%	-0.3%	-2.2%	-1.4%

The **hack** (due to A. Gill): hardcode argument cardinalities for  
build, foldr and runST.

# Compiling with optimised GHC

- We compiled GHC *itself* with cardinality optimisations;
- Then we measured improvement in *compilation runtimes*.

Program	LOC	GHC Alloc $\Delta$		GHC RT $\Delta$	
		No hack	Hack	No hack	Hack
anna	5740	-1.6%	-1.5%	-0.8%	-0.4%
cacheprof	1600	-1.7%	-0.4%	-2.3%	-1.8%
fluid	1579	-1.9%	-1.9%	-2.8%	-1.6%
gamteb	1933	-0.5%	-0.1%	-0.5%	-0.1%
parser	2379	-0.7%	-0.2%	-2.6%	-0.6%
veritas	4674	-1.4%	-0.3%	-4.5%	-4.1%

# Beyond GHC optimisations



## **Program Synthesis by Type-Guided Abstraction Refinement**

ZHENG GUO, UC San Diego, USA  
MICHAEL JAMES, UC San Diego, USA  
DAVID JUSTO, UC San Diego, USA  
JIAXIAO ZHOU, UC San Diego, USA  
ZITENG WANG, UC San Diego, USA  
RANJIT JHALA, UC San Diego, USA  
NADIA POLIKARPOVA, UC San Diego, USA

- Hoogle+ synthesis algorithm (POPL'20) relies on cardinality analysis to eliminate terms where some of the inputs are unused.

# To take away

- **Cardinality analysis** is **simple** to design and understand: it's all about *usage demands* and *demand transformers*;
- It is **cheap to implement**: we added only 250 LOC to GHC;
- It is conservative, which makes it **fast** and **modular**;
- *Call demands* make it **higher-order**, so the analysis can infer demands on higher-order function arguments;
- It is **reasonably efficient**: optimised GHC compiles up to 4% faster.
- The ideas of cardinality analysis extend beyond just optimisations in GHC.

Thanks!