

A semantics for imprecise exceptions

Simon Peyton Jones
Microsoft Research Ltd, Cambridge
simonpj@microsoft.com

Alastair Reid
Yale University
reid-alastair@cs.yale.edu

Tony Hoare
Cambridge University

Simon Marlow
Glasgow University
simonm@dcs.gla.ac.uk

October 28, 1998

Abstract

Some modern superscalar microprocessors provide only *imprecise exceptions*. That is, they do not guarantee to report the same exception that would be encountered by a straightforward sequential execution of the program. In exchange, they offer increased performance or decreased area (which amount to much the same thing).

This performance/precision tradeoff has not so far been much explored at the programming language level. In this paper we propose a design for imprecise exceptions in the lazy functional programming language Haskell. We discuss various simpler designs, and conclude that imprecision is essential if the language is still to enjoy its current rich algebra of transformations. We sketch a precise semantics for the language extended with exceptions.

From the functional programming point of view, the paper shows how to extend Haskell with exceptions without crippling the language or its compilers. From the point of view of the wider programming language community, we pose the question of whether precision and performance can be traded off in other languages too.

This paper has been submitted to PLDI'99.

1 Introduction

All current programming languages that support exceptions take it for granted that the language definition should specify, for a given program, what exception, if any, is raised when the program is executed. That used to be the case in microprocessor architecture too, but it is no longer so. Some processors, notably the Alpha, provide so-called *imprecise exceptions*. These CPUs execute many instructions in parallel, and perhaps out of order; it follows that the first exception (divide-by-zero, say) that is encountered is not necessarily the first that would be encountered in simple sequential execution. One approach is to provide lots of hardware to sort the mess out, and maintain the programmer's illusion of a simple sequential execution engine; this is what the Pentium does. Another, taken by the Alpha, is to give a less precise indication of the whereabouts of the exception.

In this paper we explore this same idea at the level of the

programming language. The compiler, or the programmer, might want to improve performance by changing the program's evaluation order. But changing the evaluation order may change which exception is encountered first. One solution is to ban such transformations, or to restrict them to evaluations that provably cannot raise exceptions. The alternative we propose here is to trade precision for performance: permit richer transformations, and make the language semantics less precise with respect to which exception is raised.

We make all this concrete by considering a particular programming language, Haskell, that currently lacks exceptions. Our contributions are as follows:

- We review and critique the folk-lore on exception-handling in a lazy language like Haskell (Section 2). Non-functional programmers may find the idea of exceptions-as-values, as opposed to exceptions-as-control-flow, interesting.
- We present a new design, based on *sets* of exceptions, to model imprecision about which exceptions can occur (Section 3).
- We sketch a semantics for the resulting language, using two layers: a denotational semantics for pure expressions (including exception-raising ones), and an operational semantics “on top” that deals with exception handling, as well as input/output (Section 4).
- Informed by this semantics, we show that various extensions of the basic idea, such as resource-exhaustion interrupts, can readily be accommodated; and that others, such as a “pure” exception handler, cannot (Section 5).

There has been a small flurry of recent proposals and papers on exception-handling in Haskell [2, 10, 9]. The distinctive feature of this paper is its focus on the semantics of the resulting language. The trick lies in getting the nice features of exceptions (efficiency, implicit propagation, and the like) without throwing the baby out with the bath water and crippling the language design.

Those less interested in functional programming *per se* may nevertheless find interesting our development of the (old) idea of exceptions-as-values, and the trade-off between precision and performance.

2 The status quo ante

Haskell has managed without exceptions for a long time, so it is natural to ask whether they are either necessary or appropriate. We briefly explore this question, as a way of setting the scene for the rest of the paper.

Before we begin, it is worth identifying three different ways in which exceptions are typically used in languages that support them:

Disaster recovery uses an exception to signal a (hopefully rare) error condition, such as division by zero or an assertion failure. In a language like ML or Haskell we may add pattern-match failure, when a function is applied to a value for which it does not have a defining equation (e.g. `head` of the empty list). The programmer can usually also raise an exception, using a primitive such as `raise`.

The exception handler typically catches exceptions from a large chunk of code, and performs some kind of recovery action.

Exception handling used in this way provides a degree of modularity: one part of a system can protect itself against failure in another part of the system.

Alternative return. Exceptions are sometimes used as an alternative way to return a value from a function, where no error condition is necessarily implied. An example might be looking up a key in a finite map: it's not necessarily an error if the key isn't in the map, but in languages that support exceptions it's not unusual to see them used in this way.

The exception handler typically catches exceptions from a relatively circumscribed chunk of code, and serves mainly as an alternative continuation for a call.

Asynchronous events. In some languages, an asynchronous external event, such as the programmer typing “`C`” or a timeout, are reflected into the programmer's model as an exception. We call such things *asynchronous exceptions*, to distinguish them from the two previous categories, which are both *synchronous exceptions*.

2.1 Exceptions as values

No lazy functional programming language has so far supported exceptions, for two apparently persuasive reasons.

Firstly, lazy evaluation scrambles control flow. In a so-called “lazy” functional language, evaluation is demand-driven; that is, an expression is evaluated only when its value is required [11]. As a result, programs don't have a readily-predictable control flow; the only productive way to think about an expression is to consider the value it computes, not the way in which it computes it. Since exceptions are typically explained in terms of changes in control flow, exceptions and lazy evaluation do not appear very compatible.

Secondly, exceptions can be explicitly encoded in values, in the existing language, so perhaps exceptions are in any case unnecessary. For example, consider a function, `f`, that takes

an integer argument, and either returns an integer or raises an exception. We can encode it in Haskell thus:

```
data ExVal a = OK a
             | Bad Exception

f :: Int -> ExVal Int
f x = ...defn of f...
```

The `data` declaration says that a value of type `ExVal t` is either of the form `(Bad ex)`, where `ex` has type `Exception`, or is of the form `(OK val)`, where `val` has type `t`. That is, *the exception is encoded into the value*. The type signature of `f` declares that `f` returns a result of type `ExVal Int`; that is, either an `Int` or an exception value.

Any consumer of `f`'s result is forced, willy nilly, to first perform a case analysis on it:

```
case (f 3) of
  OK val -> ...normal case...
  Bad ex -> ...handle exception...
```

There are good things about this approach: no extension to the language is necessary; the type of a function makes it clear whether it can raise an exception; and the type system makes it impossible to forget to handle an exception.

The idea of exceptions as values is very old. Wadler's influential paper “How to replace failure with a list of successes” is an early example [7]. Subsequently it was realised that the exception type constructor, `ExVal`, forms a *monad* [4, 6]. Rather than having lots of *ad hoc* pattern matches on `OK` and `Bad`, standard monadic machinery such as Haskell's `do` notation, can hide away much of the plumbing.

2.2 Inadequacies of exceptions as values

Encoding exceptions explicitly in an un-modified language works beautifully for the alternative-return usage of exceptions, but badly for the disaster-recovery use, and not at all for asynchronous events. There are several distinct problems:

- *Increased strictness.* When adding exception handling to a lazy program, it is very easy to accidentally make the program strict by testing function arguments for errors when they are applied instead of when they are used. For example, parser combinators typically don't return any result until the entire input has been parsed — this may be a major problem on very large (or infinite!) input streams.

This increased strictness is often easily avoided by using the call-by-name translation instead of the call-by-value translation — but it's surprising how easy it is to fall into this trap.

- *Excessive clutter.* The principal feature of an exception mechanism is that exceptions propagate implicitly, without requiring extra clutter in the code between the place the exception is raised and where it is handled. In stark contrast, the explicit-encoding approach forces all the intermediate code to deal explicitly (or monadically) with exceptional values. The

resulting clutter is absolutely intolerable for those situations where exceptions are used to signal disaster, because in these cases propagation is almost always required. For example, where we would originally have written:

```
(f x) + (g y)
```

we are now forced to write¹:

```
case (f x) of
  Bad ex -> Bad ex
  OK xv -> case (g y) of
    Bad ex -> Bad ex
    OK yv -> OK (xv+yv)
```

These strictures do not apply where exceptions are used as an alternative return mechanism. In this case, the approach works beautifully because propagation isn't nearly so important.

- *Built-in exceptions are un-catchable.* In Haskell, all the causes of failure recognised by the language itself (such as divide by zero, and pattern-match failure) are treated semantically as bottom (\perp), and are treated in practice by bringing the program to a halt. Haskell allows the program to trigger a similar failure by calling the standard function `error`. So, evaluating the call:

```
error "Urk"
```

halts execution, printing “Urk” on standard error. The language offers no way to catch and recover from any of these (synchronous) events. This is a serious problem when writing programs composed out of large pieces over which one has little control; there is just no way to recover from failure in any sub-component.

- *Poor efficiency.* Exceptions should cost very little if they don't actually occur. Alas, an explicit encoding into Haskell values forces a test-and-propagate at every call site, with a substantial cost in code size and speed.
- *Loss of transformations.* Programs written in a monadic style have many fewer transformations than their pure counterparts. We elaborate on this problem in Section 3.
- *No asynchronous exceptions.* Asynchronous exceptions, by their nature, have nothing to do with the value of the unfortunate expression that happens to be under evaluation when the external event occurs. Since they arise from external sources, they clearly cannot be dealt with as an explicitly-encoded value.

2.3 Goals

With these thoughts in mind, we have the following goals:

- Haskell programs that don't invoke exceptions should have unchanged semantics (no clutter), and run with unchanged efficiency.

¹The monadic version is nearly as bad.

- All transformations that are valid for ordinary Haskell programs should be valid for the language extended with exceptions.
- It should be possible to reason about which exceptions a program might raise. At the very least, we should be able to show that non-recursive programs will terminate, and programs that don't use arithmetic can't raise division by zero.
- In so far as non-determinism arises, it should be possible for the programmer to confine it.

These properties may seem obvious but they are a little tricky to achieve. In existing languages that support exceptions, such as ML or Ada, the need to maintain the exception semantics noticeably constrains the valid set of transformations and optimisations that a programmer or compiler can perform. Compilers often infer the set of possible exceptions with a view to lifting these restrictions, but they must be pessimistic across module boundaries in the presence of separate compilation. We claim that our design retains almost all useful opportunities for transformation, without requiring any analysis at all.

3 A new design

Adding exceptions to a lazy language, as opposed to encoding exceptions in the un-extended language, has received relatively little attention until recently. Dornan and Hammond discussed adding exceptions to the pure (non-I/O) part of a lazy language [1], and there has been a flurry of recent activity [2, 10, 9]. Drawing on this work, we propose a programming interface, for an exceptions mechanism. This sets the scene for the core of our paper, the semantics for the resulting language.

3.1 The basic idea

As discussed in Section 2.1, our first design decision is more or less forced by the fact that Haskell is a lazy language: *exceptions are associated with data values, rather than with control flow.* This differs fundamentally from the standard approach to exceptions taken for imperative, or strict functional, languages, where exceptions are associated with control flow rather than with data flow. The one place that exceptions-as-values does show up in the imperative world is the NaNs (not-a-number) of the IEEE floating point standard, where certain bit-patterns encode exceptional values, which are propagated by the floating point operations [15].

So a value (of any type) is either a “normal” value, or it is an “exceptional” value. An “exceptional” value contains an exception, and we must say what that is. The data type `Exception` is the type of exceptions. It is a fixed, system-supplied, algebraic data type, defined something like this:

```
data Exception = DivideByZero
               | Overflow
               | UserError String
               ...
```

One could imagine a simpler type (e.g. encoding an exception as an integer, or a string), or a richer type (e.g. a user-extensible data type, such as is provided by ML), but this one is a useful compromise for this paper. Nothing we say depends on the exact choice of constructors in the data type; hence the “...”.

For each type a , the function `raise` maps an `Exception` into an exceptional value of type a :

```
raise :: Exception -> a
```

Here, immediately, we see a difference from the explicit-encoding approach. *Every* type in the language has the possibility of being an exceptional value — previously only values of type `ExVal t` had that possibility.

The previously-primitive function `error` can readily be defined using `raise`:

```
error :: String -> a
error str = raise (UserError str)
```

Next, we need to be able to catch exceptions. The obvious way to do that is to take a value, and determine whether or not it is an exceptional value²:

```
getException :: a -> ExVal a
```

We will see in Section 3.3 that there is a fundamental problem with `getException`, but the “repaired” version (Section 3.3) is not dissimilar, so it will serve as a Aunt Sally for now.

The whole point of exceptions is, of course, that they propagate automatically. So integer addition, for example, should deliver an exceptional value if either of its arguments is an exceptional value — we will return to the question of what the exact semantics of propagation might be. Needless to say, propagation only happens in a strict context. For example, consider the `zipWith` function:

```
zipWith f [] [] = []
zipWith f (x:xs) (y:ys) = f x y : zipWith f xs ys
zipWith f xs ys = error "Unequal lists"
```

A call to `zipWith` may return an exception value — for example, `zipWith (+) [] [1]`. It may also return a list with an exception value at the end — for example, `zipWith (+) [1] [1,2]`. Finally, it may deliver a list whose spine is fully defined, but some of whose elements are exceptional values — for example `zipWith (/) [1,2] [1,0]`.

To repeat: it is *values* not *calls* that may be exceptional, and exceptional values may, for example, hide inside lazy data structures. To be sure that a data structure contains no exceptional values one must evaluate all the elements of that structure, with `seq` perhaps.

3.2 Implementation

One advantage of the story so far is that it is readily, and cheaply, implementable. We certainly do not want the space and time cost of explicitly tagging every value with an indication of whether it is “normal” or “exceptional”. Fortu-

²Recall that the data type `ExVal` was defined in Section 2.1

nately, the standard exception-handling mechanisms from procedural languages work perfectly well:

- `getException` forces the evaluation of its argument to head normal form; before it begins this evaluation, it marks the evaluation stack in some way.
- `raise ex` simply trims the stack to the top-most `getException` mark, and returns `Bad ex` as the result of `getException`.
- If the evaluation of the argument to `getException` completes without provoking a call to `raise`, then `getException` returns `OK val`, where `val` is the value of the argument.

Actually, matters are not quite as simple as we suggest here. In particular, trimming the stack after a call (`raise ex`) must be careful to overwrite each thunk that is under evaluation with (`raise ex`). That way, if the thunk is evaluated again, the same exception will be raised again, which is as it should be³. The details are described by [9], and need not concern us here.

The main point is that the efficiency of programs that do not invoke exceptions is unaffected. Indeed, the efficiency of any function that does not invoke exceptions explicitly is unaffected. Notice that an exceptional value *behaves* as a first class value, but it is never *explicitly represented* as such. When an exception occurs, instead of building a value that represents it, we look for the exception handler right away. The semantic model (exceptional values) is quite different from the implementation (evaluation stacks and stack trimming). The situation is similar to that with lazy evaluation: a value may *behave* as an infinite list, but it is certainly never *explicitly represented* as such.

3.3 A problem and its solution

There is a well-known difficulty with the approach we have just described: it invalidates many useful transformations. For example, integer addition should be commutative; that is, $e_1 + e_2 = e_2 + e_1$. But what are we to make of this expression?

```
getException ((1/0) + (error "Urk"))
```

Does it deliver `DivideByZero` or `UserError "Urk"`? Urk indeed! There are two well known approaches, and one more cunning one which we shall adopt:

- Fix the evaluation order of `+`, and declare that if its first argument is exceptional then that’s the exception that is returned. This is essentially the approach of [1]. It gives rise to a simple semantics, but has the Very Bad Feature that it invalidates many useful transformations – in particular, ones that alter the order of evaluation. That in turn leads to elaborate analyses to try to identify the common case where the first argument cannot raise an exception, and so on.

³Real implementations have to overwrite a thunk with a “black hole” when its evaluation is begun to avoid a celebrated space leak [3]. That is why, when an exception causes their evaluation to be abandoned, they must be overwritten with something more informative.

- Go non-deterministic. That is, declare that `+` makes a non-deterministic choice of which argument to evaluate first. Then the compiler is free to make that choice however it likes. This approach exposes non-determinism in the source language, which also invalidates useful laws. In particular, β reduction is not valid any more. For example, consider:

```
let x = (1/0) + (error "Urk")
in getException x == getException x
```

As it stands, the value of this expression is presumably `True`. But if the two occurrences of `x` are each replaced by `x`'s right hand side, then the non-deterministic `+` might (in principle) make a different choice at its two occurrences, so the expression could be `False`. We count this too high a price to pay.

- The more cunning choice is to return *both* exceptions! That is, we redefine an exceptional value to contain a *set* of exceptions, instead of just one; and `+` takes the union of the exception sets of its two arguments. Now `(1/0) + (error "Urk")` returns an exceptional value including both `DivideByZero` and `UserError "Urk"`, and (semantically) it will do so regardless of the order in which `+` evaluates its arguments.

The allegedly cunning choice may have fixed the commutativity of `+`, but we have to worry about several other inter-related things now. For a start, what is `getException` to do? There are two possibilities.

First, `getException` could return the set of exceptions. This is a real disaster from an implementation point of view! It would mean that the implementation would really have to maintain a set of exceptions; if the first argument to `+` failed, then the second would have to be evaluated anyway so that any exceptions in it could be gathered up.

The cunning choice is only cunning because there is another alternative: `getException` can choose just one member of the set of exceptions to return. Of course, that simply exposes the non-determinism again, but we can employ a standard trick: put `getException` in the `I0` monad. Thus:

```
getException :: a -> IO (ExVal a)
```

In Haskell, a value of type `IO t` is a *computation* that might perform some input/output, before eventually returning a value of type `t`. A value of type `IO t` is a first-class value — it can be passed as an argument, stored in a data structure — and evaluating it has no side effects. Only when it is *performed* does it have an effect. An entire Haskell program is a single value of type `IO ()`; to run the program is to perform the specified computation. For example, here is a complete Haskell program that gets one character from standard input and echos it to standard output⁴:

```
main :: IO ()
main = getChar      >>= (\ch ->
      putChar ch    >>= (\() ->
      return ()
    ))
```

The types of the various functions involved are as follows:

```
(>>=)  :: IO a -> (a -> IO b) -> IO b
return :: a -> IO a
getChar :: IO Char
putChar :: Char -> IO ()
```

The combinators `>>=` glues together two `IO` computations in sequence, passing the result from the first to the second. `return` does no input/output, simply returning its argument. `getChar` gets a character from the standard input and returns it; `putChar` does the reverse. When `main` is performed, it performs `getChar`, reading a character from standard input, and then performs the computation obtained by applying the `\ch -> ...` abstraction to the character, in this case `putChar ch`. A more complete discussion of monadic I/O can be found in [14].

By giving `getException` an `IO` type we allow it (in principle) to perform input/output. Hence, when choosing which of the exceptions in the set to choose, `getException` is free to consult some external oracle (the FT Share Index, say). Each call to `getException` can make a different choice; the same call to `getException` in different runs of the same program can make a different choice; and so on.

Beta reduction remains valid. For example the meaning of:

```
let
  x = (1/0) + error "Urk"
in
getException x >>= (\v1 ->
getException x >>= (\v2 ->
return (v1==v2)))
```

is unaffected if both occurrences of `x` are replaced by `x`'s right hand side, thus:

```
getException ((1/0) + error "Urk") >>= (\v1 ->
getException ((1/0) + error "Urk") >>= (\v2 ->
return (v1==v2)))
```

Why? Because whether or not this substitution is made, `getException` will be performed twice, making an independent non-deterministic choice each time. Like any `IO` computation, (`getException e`) can be evaluated and shared without actually performing the I/O itself. That only happens when the computation is performed.

The really nice thing about this approach is that the stack-trimming implementation does not have to change. The set of exceptions associated with an exceptional value is represented by a single member, namely the exception that happens to be encountered first. `getException` works just as before: mark the evaluation stack, and evaluate its argument. Successive runs of a program, using the same compiler optimisation level, will in practice give the same behaviour; but if the program is recompiled with different optimisation settings, then indeed the order of evaluation might change, so a different exception might be encountered first, and so the exception returned by `getException` might change.

This idea is based on an old paper by Hughes and O'Donnell [12], but it was Fergus Henderson who first suggested how it could be adapted to the setting of exception handling [2]. He was the first to realise that non-determinism in the *exceptions* could be kept separate from non-determinism in the normal *values* of a program.

⁴The “\” is Haskell's notation for λ

4 Semantics

Everything so far in this paper is either well-known or folklore. Our main contribution is to give a semantics to these non-deterministic exceptions, which we do in this section. Here are two difficulties.

- Consider

```
undefined + error "Urk"
```

Here, `undefined` is any expression whose evaluation diverges. It might be declared like this:

```
undefined = f True
  where
    f x = f (not x)
```

So, does `(undefined + error "Urk")` loop forever, or does it return an exceptional value? Answer: it all depends on the evaluation order of `+`. As is often the case, bottom muddies the waters.

- Is the following equation true?

```
case x of
  (a,b) -> case y of
    (p,q) -> e
=
case y of
  (p,q) -> case x of
    (a,b) -> e
```

In Haskell the answer is “yes”; since we are going to evaluate both `x` and `y`, it doesn’t matter which order we evaluate them in. Indeed, the whole point of strictness analysis is to figure out which things are sure to be evaluated in the end, so that they can be evaluated in advance [5]. But if `x` and `y` are both bound to exceptional values, then the order of the `cases` clearly determines which exception will be encountered. Unlike the `+` case, it is far from obvious how to combine the exceptional value sets for `x` and `y`: in general the right hand side of a `case` alternative might depend on the variables bound in the pattern, and it would be unpleasant for the semantics to depend on that.

The rest of this section gives a denotational semantics for Haskell extended with exceptions, that addresses both of these problems. We solve the first by identifying \perp with the set of all possible exceptions; we solve the latter by (semantically) evaluating the `case` alternatives in “exception-finding mode”.

4.1 Domains

First we describe the domain $\llbracket \tau \rrbracket$ that is associated with each Haskell type τ . We use a rather standard monadic translation, for a monad M , defined thus:

```
 $\mathcal{M} t = t_{\perp} + \mathcal{P}(\mathcal{E})_{\perp}$ 
 $\mathcal{E} = \{\text{DivideByZero, Overflow, UserError, ...}\}$ 
```

The “+” in this equation is coalesced sum; that is, the bottom element of $\llbracket \tau \rrbracket_{\perp}$ is coalesced with the bottom element of $\mathcal{P}(\mathcal{E})_{\perp}$. The set \mathcal{E} is the set of all the possible synchronous exceptions; to simplify the semantics we neglect the `String` argument to `UserError`. $\mathcal{P}(\mathcal{E})$ is the lattice of all subsets of \mathcal{E} , under the ordering

$$s_1 \sqsubseteq s_2 \equiv s_1 \supseteq s_2$$

That is, the bottom element is the set \mathcal{E} , and the top element is the empty set. This corresponds to the idea that the fewer exceptions that are in the exceptional value, the more information the value contains. The least informative value contains all exceptions. This entire lattice is lifted, by adding an extra bottom element, which we also identify with a set of exceptions:

$$\perp = \mathcal{E} \cup \{\text{NonTermination}\}$$

At first we distinguished \perp from the set of all exceptions, but that turns out not to work. Instead, we identify them, adding one new constructor, `NonTermination`, to the `Exception` type:

```
data Exception = ... -- (as before)
               | NonTermination
```

This construction of $\mathcal{P}(\mathcal{E})_{\perp}$ is a very standard semantic coding trick; it is closely analogous to a canonical representation of the Smyth powerdomain over a flat domain, given by [8].

Here is a more element-wise way to define \mathcal{M} , in which we tag “normal” values with `Ok`, and “exceptional” values (including \perp) with `Bad`:

$$\begin{aligned} \mathcal{M} t &= \{Ok\ v \mid v \in t\} \cup \\ &\quad \{Bad\ s \mid s \subseteq \mathcal{E}\} \cup \\ &\quad \{Bad\ (\mathcal{E} \cup \{\text{NonTermination}\})\} \end{aligned}$$

One might wonder what sort of a value `Bad {}` is: what is an exceptional value containing the empty set of exceptions? Indeed, such a value cannot be the denotation of any term, but we will see shortly that it is a very useful value nevertheless.

Now that we have constructed the exception monad, we can translate Haskell types into domains in the usual way:

$$\begin{aligned} \llbracket \text{Int} \rrbracket &= \mathcal{M} \mathcal{Z} \\ \llbracket \tau_1 \rightarrow \tau_2 \rrbracket &= \mathcal{M} (\llbracket \tau_1 \rrbracket \rightarrow \llbracket \tau_2 \rrbracket) \\ \llbracket \langle \tau_1, \tau_2 \rangle \rrbracket &= \mathcal{M} (\llbracket \tau_1 \rrbracket \times \llbracket \tau_2 \rrbracket) \\ &\dots \text{etc} \dots \end{aligned}$$

We refrain from giving the complete encoding for arbitrary recursive data types, which is complicated. The point is that we simply replace the normal Haskell monad, namely lifting, with our new monad M .

4.2 Combinators

Next, we must give the denotation, or meaning, of each form of language expression. Figure 1 gives the syntax of the small language we treat here. The denotation of an expression e in an environment ρ is written $\llbracket e \rrbracket \rho$.

$e ::= x$	variable
k	constant
$e_1 e_2$	application
$\lambda x.e$	abstraction
$C e_1 \dots e_n$	constructors
$\text{case } e \text{ of } \{ \dots p_i \rightarrow r_i; \dots \}$	matching
$\text{raise } C$	raise exception
$e_1 + e_2$	primitives
$\text{fix } e$	fixpoint
$p ::= C x_1 \dots x_n$	pattern

Figure 1: Syntax of a tiny language

We start with $+$:

$$\begin{aligned} \llbracket e_1 + e_2 \rrbracket \rho &= v_1 \oplus v_2 && \text{if } Ok\ v_1 = \llbracket e_1 \rrbracket \rho \\ & && \text{and } Ok\ v_2 = \llbracket e_2 \rrbracket \rho \\ &= Bad\ (\mathcal{S}(\llbracket e_1 \rrbracket \rho) \cup \mathcal{S}(\llbracket e_2 \rrbracket \rho)) && \text{otherwise} \end{aligned}$$

The first equation is used if both arguments are normal values. The second is used if either argument is an exceptional value, in which case the exceptions from the two arguments are unioned. We use the auxiliary function $\mathcal{S}()$, which returns the empty set for a normal value, and the set of exceptions for an exceptional value:

$$\begin{aligned} \mathcal{S}(Ok\ v) &= \emptyset \\ \mathcal{S}(Bad\ s) &= s \end{aligned}$$

The auxiliary function \oplus simply does addition, checking for overflow:

$$\begin{aligned} v_1 \oplus v_2 &= Ok\ (v_1 + v_2) && \text{if } \perp 2^{31} \langle (v_1 + v_2) \rangle 2^{31} \\ &= Bad\ \{\text{0verflow}\} && \text{otherwise} \end{aligned}$$

The definition of $\llbracket + \rrbracket$ is monotonic with respect to \sqsubseteq , as it must be. The fact that $+$ is strict in both arguments is a consequence of the fact that \perp is the set of all exceptions; a moment's thought should convince you that if either argument is this set then so is the result.

Next, we deal with **raise**:

$$\llbracket \text{raise } C \rrbracket \rho = Bad\ \{C\}$$

Now we can understand the semantics of the problematic expression given above:

```
undefined + error "Urk"
```

Its meaning is the union of the set of all exceptions (which is the value of **undefined**), and the singleton set **UserError "Urk"**, which is of course just \perp , the set of all exceptions.

The rules for function abstraction and application are interesting:

$$\begin{aligned} \llbracket \lambda x.e \rrbracket \rho &= Ok\ (\lambda y.\llbracket e \rrbracket \rho[y/x]) \\ \llbracket e_1 e_2 \rrbracket \rho &= f\ (\llbracket e_2 \rrbracket \rho) && \text{if } Ok\ f = \llbracket e_1 \rrbracket \rho \\ &= Bad\ (s \cup \mathcal{S}(\llbracket e_2 \rrbracket \rho)) && \text{if } Bad\ s = \llbracket e_1 \rrbracket \rho \end{aligned}$$

A lambda abstraction is a normal value; that is $\lambda x.\perp \neq \perp$. The (more purist) identification of these two values is

impossible implement: how can **getException** distinguish $\lambda x.\perp$ from $\lambda x.v$, where $v \neq \perp$? Fortunately, in Haskell $\lambda x.\perp$ and \perp are already distinct.

Applying normal function to a value is straightforward, but matters are more interesting if the function is an exceptional value. In this case we must union its exception set with that of its argument, because under some circumstances (notably if the function is strict) we might legitimately evaluate the argument first; if we neglected to union in the argument's exceptions, the semantics would not allow this change to call by value. Notice that we must *not* union in the argument's exceptions if the function is a normal value, or else we would lose β reduction; consider $(\lambda x.3)(1/0)$

The rules for constants and constructor applications are simple; they both return normal values. Constructors are non-strict, and hence do not propagate exceptions in their arguments. Variables and fixpoints are also easy.

$$\begin{aligned} \llbracket k \rrbracket \rho &= Ok\ k \\ \llbracket C\ e_1 \dots e_n \rrbracket \rho &= Ok\ (C(\llbracket e_1 \rrbracket \rho, \dots, \llbracket e_n \rrbracket \rho)) \\ \llbracket x \rrbracket \rho &= \rho(x) \\ \llbracket \text{fix } e \rrbracket \rho &= \bigsqcup_{k=0}^{\infty} (\llbracket e \rrbracket \rho)^k(\perp) \end{aligned}$$

4.3 case expressions

Lastly we turn our attention to **case** expressions. In fact, in the Glasgow Haskell Compiler (GHC), $+$ is implemented in terms of **case** expressions, thus:

```
data Int = I# Int#

a + b = case a of
  I# v -> case b of
    I# w -> I# (v +# w)
```

In GHC, **Int** is an algebraic data type with a single constructor **I#**, whose argument is an unboxed integer of type **Int#**. We had better make sure that **case** expressions are dealt with in a way that is consistent with the direct semantics for $+$ given earlier.

Here, then, is the slightly surprising rule for **case**:

$$\begin{aligned} \llbracket \text{case } e \text{ of } \{ p_i \rightarrow r_i \} \rrbracket \rho &= \llbracket r_i \rrbracket \rho[v/p_i] \\ &\quad \text{if } Ok\ v = \llbracket e \rrbracket \rho \\ &\quad \text{and } v \text{ matches } p_i \\ &= Bad\ (s \cup \bigcup_i \mathcal{S}(\llbracket r_i \rrbracket \rho[Bad\ \{ \} / p_i])) \\ &\quad \text{if } Bad\ s = \llbracket e \rrbracket \rho \end{aligned}$$

The first case is the usual one: if the case scrutinee evaluates to a "normal" value v , then select the appropriate case alternative. The notation is a little informal: $\rho[v/p_i]$ means the environment ρ with the free variables of the pattern p_i bound to the appropriate components of v .

The second equation is the interesting one. If the scrutinee turns out to be a set of exceptions (which, recall, includes \perp), the obvious thing to do is to return just that set — but

doing so would invalidate the case-switching transformation. Intuitively, we want to evaluate all the branches anyway, in “exception-finding mode”. We model this by taking the denotations of all the right hand sides, binding each of the pattern-bound variables to the strange value *Bad* $\{\}$. Then we union together all the exception sets that result, along with the exception set from the scrutinee. The idea is exactly the same as in the special case of $+$, and function application: if the first argument of $+$ raises an exception we still union in the exceptions from the second argument. Here, if the case scrutinee raises an exception, we still union in the exceptions from the alternatives.

Remember that there is no implication that an implementation will do anything other than return the first exception that happens to be encountered. The rather curious semantics is necessary, though, to validate transformations that change the order of evaluation, such as that given earlier:

```

case x of
  (a,b) -> case y of
            (p,q) -> e
=
case y of
  (p,q) -> case x of
            (a,b) -> e

```

4.4 Semantics of `getException`

So far we have not mentioned `getException`. The semantics of operations in the `I0` monad, such as `getException`, necessarily involve input/output and non-determinism. The most straightforward way of modelling these aspects is by giving an *operational* semantics for the `I0` layer, in contrast to the *denotational* semantics we have given for the purely-functional layer.

We give the operational semantics as follows. From a semantic point of view we regard `I0` as an algebraic data type with constructors `return`, `>>=`, `putChar`, `getChar`, `getException`. The behaviour of a program is the set of traces obtained from the following labelled transition system, which acts on the *denotation* of the program. One advantage of this presentation is that it scales to other extensions, such as adding concurrency to the language [13].

Here are the structural transition rules:

$$\frac{v_1 \rightarrow v_2}{(v_1 \gg= k) \rightarrow (v_2 \gg= k)}$$

$$((\text{return } v) \gg= k) \rightarrow (k v)$$

The first ensures that transitions can occur inside the first operand of the `>>=` constructor; the second explains that a `return` constructor just passes its value to the second argument of the enclosing `>>=`. The rules for input/output are now quite simple:

```

getChar   $\xrightarrow{?c}$  return c
putChar c  $\xrightarrow{!c}$  return ()

```

The “ $?c$ ” on top of the arrow indicates that the transition takes place by reading a character c from the environment; and inversely for “ $!c$ ”.

Now we can get to the semantics of exceptions. The rules are:

```

getException (Ok v)  → return (OK (Ok v))
getException (Bad s) → return (Bad x)
                        if x ∈ s
getException (Bad s) → getException (Bad s)
                        if NonTermination ∈ s

```

If `getException` scrutinises a “normal” value, it just returns it, wrapped in an `OK` constructor.

For “exceptional” values, there are two choices: either

- pick an arbitrary member of the set of exceptions and return it, or
- if `NonTermination` is in the set of exceptions, then make a transition to the same state.

The transition rules for `getException` are deliberately non-deterministic. In particular, if the argument to `getException` is \perp , then `getException` may diverge, or it may return an arbitrary exception.

5 Variations on the theme

5.1 Pure functions on exceptional values

Is it possible to do anything with an exceptional value other than choose an exception from it with `getException`? Following [12], one possibility suggests itself as a new primitive function (i.e. one not definable with the primitives so far described):

```
mapException :: (Exception -> Exception) -> a -> a
```

Semantically, `mapException` applies its functional argument to each member of the set of exceptions (if any) in its second argument; it does nothing to normal values. From an implementation point of view, it applies the function to the sole representative (if any) of that set. Here’s an example of using `mapException` to catch all exceptions in `e` and raise `UserError "Urk"` instead:

```
mapException (\e -> UserError "Urk") e
```

Notice that `mapException` does not need to be in the `I0` monad to preserve determinism.

`mapException` maps one kind of exception to another, but it doesn’t let us get from exceptions back into normal values. Is it possible to go further? Is it possible, for example, to ask “is this an exceptional value”?

```
isException :: a -> Bool
```

(It would be easy to define `isException` with a monadic type `a -> I0 Bool`; the question is whether it can have a pure, non-monadic, type.) At first `isException` looks reasonable, because it hides just which exception is being raised — but it doesn’t work. What is the value of the following expression?

```
isException ((1/0) + undefined)
```

Since the compiler might choose to evaluate the second argument of `+` first, the result had better be \perp ; but if the compiler evaluates the first argument of `+` first, the result will be `True`. In short, `isException` seems to be reasonable semantically (e.g. we can give it a monotonic definition), but it is not implementable⁵. A variant of exactly this problem is raised by Reid [9]; the solution is not to provide `isException`.

5.2 Asynchronous exceptions

All the exceptions we have discussed so far are synchronous exceptions (Section 2). If the evaluation of an expression yields a synchronous exception, then another evaluation of the same expression will yield the same result. But what about asynchronous exceptions, such as interrupts and resource-related failures (e.g. timeout, stack overflow, heap exhaustion)? They differ from synchronous exceptions in that they may *not* recur (at all) if the same program is run again. It is obviously inappropriate to regard such exceptions as part of the denotation of an expression.

Fortunately, they can fit in the same general framework. We have to enrich the `Exception` type with constructors indicating the cause of the exception. Then we simply add to `getException`'s abilities. Since `getException` is in the `IO` monad, it can easily say “if the evaluation of my argument goes on for too long, I will terminate evaluation and return `Bad Timeout`”, and similarly for interrupts and so on.

```
getException v → return (Bad x)
                if x is an asynchronous exception
```

This rule says that, v , the argument of `getException`, may be ignored. v might not be an exceptional value — it might be say, `42` — but `getException` is free to discard it and return any asynchronous exception instead. In the case of timeout, `getException` would perform this transition if v took too long to compute; in the case of a keyboard interrupt, `getException` would discard v immediately; and so on.

There is a fascinating wrinkle in the implementation of external exceptions: when trimming the stack, we must overwrite each thunk under evaluation with a kind of “resumable continuation”, rather than a computation which raises the exception again. The details are in [10].

The possibility of asynchronous interrupts is another reason to be suspicious of `isException`.

5.3 Detectable bottoms

There are some sorts of divergence that are detectable by a compiler or its runtime system. For example, suppose that `undefined` was declared like this:

```
black = black + 1
```

Here, `black` is readily detected as a so-called “black hole” by many graph reduction implementations. Under these circumstances, `getException black` is permitted, but not required, to return `Bad NonTermination` instead of going into

⁵The situation reminds us of the celebrated “parallel or” operator.

a loop! Whether or not it does so is an implementation choice — perhaps implementations will compete on the skill with which they detect such errors.

5.4 Non-deterministic sets

Henderson's original proposal called for a separate data type for *non-empty, non-deterministic sets*, à la Hughes and O'Donnell, with the exception library making use of this data type. In particular, the operation of choosing a member of a non-deterministic set is the one that is proposed for the `IO` monad:

```
choose :: NDSet a -> IO a
```

He proposes a separate operation to extract the exceptions from a value:

```
getExceptions :: a -> Either (NDSet Exception) a
```

Recall that the `Either` type is defined like this:

```
data Either a b = Left a | Right b
```

There are at least two difficulties with this proposal:

- `getExceptions` falls foul of exactly the same problems as `isException` (Section 5.1). It is not implementable; or equivalently, what is implementable is indescribable (literally: we don't know how to describe its semantics).
- It is not clear what `getExceptions` should do for asynchronous exceptions.

In short, the way we have dealt with \perp in our semantics, which is key to our approach, pretty much precludes Henderson's proposed separation. (One could put `getExceptions` in the `IO` monad too, but then nothing has been gained over `getException`.)

6 Status and future work

6.1 Transformations

Our overall goal is to add exceptions to Haskell without losing useful transformations. Yet it cannot be true that we lose *no* transformations. For example, in Haskell as it stands, the following equation holds:

```
error "This" = error "That"
```

Why? Because both are semantically equal to \perp . In our semantics this equality no longer holds — and rightly so! So our semantics correctly distinguishes some expressions that Haskell currently identifies.

Some transformations that are identities in Haskell become refinements in our new system. Consider:

```
lhs = (case e of { True -> f; False -> g }) x
rhs = case e of { True -> (f x); False -> (g x) }
```

Using $e = \text{raise } E$, $x = \text{raise } X$, and $f = g = \lambda v.1$, we get $\llbracket lhs \rrbracket \rho = \text{Bad } \{E, X\}$ but $\llbracket rhs \rrbracket \rho = \text{Bad } \{E\}$. Hence, $lhs \sqsubseteq rhs$, but not $lhs = rhs$. We argue that it is legitimate to perform a transformation that increases information – in this case, reduces uncertainty about which exceptions can be raised.

We currently lack a systematic way to say which identities continue to hold, which turn into refinements, and which no longer hold. We conjecture that the lost laws deserve to be lost, and that optimising transformations are either identities or refinements. It would be interesting to try to formalise and prove this conjecture.

6.2 A spectrum of precision

There is actually a continuum between our semantics and the “fixed evaluation order” semantics, which fully determines which exception is raised. As one moves along the spectrum towards our proposal, more compiler transformations become valid — but there is a price to pay. That price is that the semantics becomes vaguer about which exceptions can be raised. Our view is that we should optimise for the no-exception case. In that case we want as many transformations to be valid as possible. The cost is that if something does go wrong in the program, then the semantics does not guarantee very precisely what exception will show up. An extreme case is this:

```
getException (undefined + error "Urk")
```

Since `undefined` has value \perp , `getException` is quite justified in returning `Bad Overflow`, or some other quite fictitious exception — and in principle a compiler refinement might do the same. It isn't clear whether this is a bug or a feature.

6.3 Status

As usual, implementation is ahead of theory: the Glasgow Haskell Compiler (4.0 and later) implements `raise` and `getException` just as described above. If nothing else, this reassures us that there are no hidden implementation traps. If this paper makes a contribution, we believe that it is in the area of guidance about what is, and what is not, semantically justifiable in the programming interface. For example, GHC originally implemented a version of `isException`, which we now believe to be quite wrong-headed.

Incidentally, exceptions in the `IO` monad itself are also handled in the same way, which makes the implementation of the `IO` monad very much more efficient, and very much less greedy on code space. Previously, every `>>=` operation had to test for, and propagate, exceptions.

7 Acknowledgement

We gratefully acknowledge helpful feedback from Corin Pitcher and Nick Benton.

References

- [1] C. Dornan and K. Hammond. Exception handling in lazy functional languages. Technical Report CS90/R5, Department of Computing Science, University of Glasgow, Jan 1990.
- [2] F. Henderson. Non-deterministic exceptions. Electronic mail to the `haskell` mailing list, June 1998.
- [3] R. Jones. Tail recursion without space leaks. *Journal of Functional Programming*, 2(1):73–80, Jan. 1992.
- [4] E. Moggi. Computational lambda calculus and monads. In *IEEE Symposium on Logic in Computer Science*, June 1989.
- [5] A. Mycroft. The theory and practice of transforming call-by-need into call-by-value. In *Proc 4th International Symposium on Programming*, pages 269–281. Springer Verlag LNCS 83, 1981.
- [6] PL Wadler. Comprehending monads. *Mathematical Structures in Computer Science*, 2:461–493, 1992.
- [7] PL Wadler. How to replace failure by a list of successes. In *Proc Functional Programming Languages and Computer Architecture, La Jolla*. ACM, June 1995.
- [8] G. Plotkin. Domains. Technical report, Department of Computer Science, University of Edinburgh, 1983.
- [9] A. Reid. Handling exceptions in Haskell. In *submitted to Practical Applications of Declarative Languages (PADL'99)*, 1998.
- [10] A. Reid. Putting the spine back in the Spineless Tagless G-machine: an implementation of resumable black holes. In *Proc Implementing Functional Languages Workshop 1998 (IFL'98)*. Springer Verlag LNCS (to appear 1999), Sept 1998.
- [11] RJM Hughes. Why functional programming matters. *Computer Journal*, 32(2):98–107, April 1989.
- [12] RJM Hughes and JT O'Donnell. Expressing and reasoning about non-deterministic functional programs. In K. Davis and R. Hughes, editors, *Glasgow Functional Programming Workshop*, pages 308–328. Springer Workshops in Computing, 1989.
- [13] SL Peyton Jones, AJ Gordon, and SO Finne. Concurrent Haskell. In *23rd ACM Symposium on Principles of Programming Languages, St Petersburg Beach, Florida*, pages 295–308. ACM, Jan 1996.
- [14] SL Peyton Jones and PL Wadler. Imperative functional programming. In *20th ACM Symposium on Principles of Programming Languages (POPL'93), Charleston*, pages 71–84. ACM, Jan 1993.
- [15] WJ Cody *et al.* A proposed radix- and word-length independent standard for floating point arithmetic. *IEEE Micro*, 4(4):86–100, Aug. 1984.