

Zipper-based File/OS¹

An extra-program demo at the Haskell Workshop 2005

We present a file server/OS where threading and exceptions are all realized via delimited continuations.

There are no unsafe operations, no GHC let alone Unix threads, no concurrency problems. Our threads cannot even do *IO* and cannot mutate any global state — and the type system sees to it.

¹I am very grateful to the Program Committee of the Haskell Workshop 2005 and its chair, Daan Leijen, for the opportunity to present ZipperFS as an extra talk. Many thanks to Chung-chieh Shan for discussions and invaluable help with the presentation.

Getting the first impression

- ▶ Load *ZFS.hs* into *GHCi*
- ▶ Start up the system: *main* at the *GHCi* prompt
- ▶ From some terminal: *telnet localhost 1503*
 - ▶ *ls*
 - ▶ *cat fl1*
 - ▶ *cd d1*
 - ▶ *ls*
 - ▶ *ls d11*
 - ▶ *ls d11 / d111*

That was an empty directory. This all looked like UnixFS.
However, there are no *.* and *..*
 - ▶ *ls ../ d2* – another empty dir
 - ▶ *cat ../ d2 ../ fl2*

Absolute paths work too.

Filesystem term

```
type FileName = String
```

```
type FileCont = String
```

```
data Term = File String | Folder (Map.Map FileName Term)
```

```
  deriving Eq
```

```
data Path = Down | DownTo FileName | DownToN Int | Up | Next
```

```
  deriving (Eq, Show)
```

```
traverse :: (Monad m) =>
```

```
  (Path -> Term -> m (Maybe Term, Path)) -> Term -> m Term
```

Down is the same as *DownToN* 0 – descend to the first child

The ‘file system’ we have just seen is a zipper over a recursive data type based on *Data.Map* of this structure: *Term*. *Path* defines movement from one subterm to another.

The user needs to define a term traversal function of this signature. This is basically *mapM* over the term, in an *arbitrary* monad. It is *not* a fold — merely a map. It pays to be careful and define the function that maximally preserves sharing: see *ZipperM.hs*.

Generic Zipper

```
data DZipper r m term dir =  
  DZipper{  
    dz_dir   :: dir,  
    dz_path  :: [dir],  
    dz_term  :: term,  
    dz_k     :: CC r m (Maybe term, dir) →  
              CC r m (DZipper r m term dir)  
  }  
  | DZipDone term
```

Once we have defined *traverse* (see *ZipperM.hs*), we can make a zipper. A zipper is an *update* cursor into an immutable data structure — *term*. It is generic — note the polymorphism of *DZipper* over any *term* type and any direction data type. We don't care what *term* is and how to specify the movement from one node to another. Unlike Huet's Zipper, our zipper is independent of the data structure and depends solely on the traversal function *interface*.

Creating generic zipper

```
data HPReq r m dir =  
  HPReq dir (CC r m [dir] → CC r m (HPReq r m dir))  
  
dzip' term term = do  
  p ← newPrompt; path_pr ← newPrompt  
  pushPrompt p (acc_path [] (pushPrompt path_pr (  
    traverse (tf p path_pr) term >>= done p))) where  
  tf p path_pr dir term = do  
    path ← shiftP path_pr (λk → return (HPReq dir k))  
    shiftP p (λk → return (DZipper dir path term k))  
  acc_path path body = do  
    HPReq dir k ← body  
    let new_path = if dir ≡ Up then tail path else dir : path  
    acc_path new_path (k (return new_path))  
  done p term = abortP p (return $ DZipDone term)
```

We create the zipper via the following generic function. The function is quite short and fits on one slide.

Creating generic zipper (cont)

The function *dzip' term* relies on the delimited continuation operators from the *CC* monad transformer library by Dybvig, Peyton-Jones and Sabry. Not surprisingly, because the zipper is the manifestation of a delimited continuation reified as a *DZipper* record. The zipper maintains the path to the current location in the term. Again, we do so generically, regardless of the term.

Built-in traversal

- ▶ *cd / d2*
- ▶ *next*
a few times. Watch for the changes in the “shell prompt”
- ▶ When in File, one can do *ls*: indeed, one can *cd* into a file.
cat is the same as *ls*: both list directories and files.
- ▶ a few more *next*
when the traversal is finished, we are stuck at the root

Unlike UnixFS, our file system has a built-in traversal facility: from *each* node, we can get to the next. Furthermore, our traversal can start from any arbitrary node in the tree.

Multi-threading

- ▶ From another terminal: *telnet localhost 1503*
- ▶ Enter at the command prompt: *ls, cd d1, ls*
- ▶ Enter *ls* in the first terminal window

We have what looks like multi-threading. However, the whole server is a single Unix process, a single Unix thread and a single GHC thread.

We do not use handles; rather, we read/write sockets directly and rely on *select*.

Our “file server” is an OS, complete with the main *osloop*, “interrupt” handler and the syscall interface.

We use delimited continuations to implement our processes.

Transactional semantics

- ▶ From the first terminal
 - ▶ `cd / d2`
 - ▶ `touch nf`
 - ▶ `ls`
 - ▶ `echo "new content" > ../d2 / n2`
Error-check does work...
 - ▶ `echo "new content" > ../d2 / nf`
 - ▶ `cat nf`
 - ▶ `rm/`
 - ▶ `rm ../ d2` – can't remove itself or its own parent
 - ▶ `rm ../ d1`
 - ▶ `cd ..`
 - ▶ `ls`
Indeed, *d1* is gone.
- ▶ From the second terminal (the current directory was *d1*)
 - ▶ `ls`
 - ▶ `ls/`
Directory *d1* still exists

Transactional semantics (cont)

- ▶ From the first terminal
 - ▶ *commit*
- ▶ From the second terminal
 - ▶ *ls*
 - ▶ *ls/*
d1 still exists. If we open the third terminal, we find that *d1* is gone
 - ▶ *refresh*
 - ▶ *ls*
And now *d1* is gone
 - ▶ *rm / d2*
we can remove whole directory trees. Oops...
 - ▶ *ls*
it is gone indeed
 - ▶ *refresh*
but we can easily undo that
 - ▶ *ls*
we see that *d2* is back

Transactional semantics (cont)

- ▶ Strongest, “repeatable read” isolation mode
- ▶ Undo
- ▶ Multiple undo and snapshots are possible

We get this all for free, without any extra programming, courtesy of the zipper

Process isolation

$$\begin{aligned} \text{run}'\text{process} &:: (\forall m. \text{Monad } m \Rightarrow (\text{Prompt } r (\text{OSReq } r m)) \rightarrow \\ &\quad \text{CC } r m (\text{OSReq } r m)) \rightarrow \\ &\quad \text{Prompt } r (\text{OSReq } r IO) \rightarrow \text{CC } r IO (\text{OSReq } r IO) \\ \text{run}'\text{process } \text{body } p &= \text{pushPrompt } p (\text{body } p) \end{aligned}$$

Here is the function to run our “process”. The process function, the first argument, does not have the *IO* type.

The base monad type *m* is left polymorphic. Although a process runs eventually in the *IO* monad, the process *cannot* know that and hence cannot do any *IO* action. It must ask the “OS” by sending an *OSReq*. That means, a process function cannot mutate the *World* or any global state, and *the type system checks that!* Because processes cannot interfere with each other and with the OS, there is no need for any thread synchronization, locking, etc. We get the transactional semantics for free.

Processes and OS requests

```
data OSReq r m
  = OSRDone
  | OSRRead (ReadK r m)
  | OSRWrite String (UnitK r m)
  | OSRTrace String (UnitK r m) -- so a process can syslog
  | OSRCommit Term (UnitK r m)
  | OSRefresh (CC r m (FSZipper r m) → CC r m (OSReq r m))

type UnitK r m = CC r m () → CC r m (OSReq r m)
type ReadK r m = CC r m String → CC r m (OSReq r m)

svc p req = shiftP p (return ∘ req)
```

Here's the type of *syscalls*: read, write, commit the changes to the file system, refresh, write to the syslog.

The function *svc* is the “supervisor call”. A process invokes *svc* to request a service from the “kernel”.

Conflict resolution

Since different processes manipulate their own (copy-on-write) terms (i.e., file systems), when processes commit, there may arise conflicts.

One has to implement some conflict resolution — be it versioning, patching, asking for permission for update, etc. In our system, these policies are implemented at the level of the supervisor rather than at the level of a process. Because processes are “pure”, always ask the supervisor for anything, and the supervisor has the view of the global state, the resolution policies become easier to implement.

Smart sharing

- ▶ *quit* on all terminals, ↑C in the *GHCi* window
- ▶ At the *GHCi* prompt: *main' fs2*
- ▶ From some terminal: *telnet localhost 1503*
 - ▶ *ls*
 - ▶ *cd d1*
 - ▶ *ls*
 - ▶ *cd d1*
 - ▶ *ls*

fs2 :: Term =

Folder \$ Map.fromList [("d1", fs2), ("f11", File "File1")]

Here is the ‘file system’ *fs2*. It has a cycle: whenever we descend into *d1*, we get back into it. No surprises here: you can do that in Unix, if you are root: hard directory link.

Smart sharing (cont)

- ▶ *touch newfile*

We are in the directory */d1 / d1/*

- ▶ *cd/*

- ▶ *ls*

No *newfile* here!

- ▶ *cd d1*

- ▶ *ls*

newfile is not here either

- ▶ *cd d1*

- ▶ *ls*

Now, it is here

- ▶ *cd d1*

- ▶ *ls*

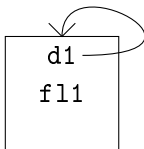
- ▶ *ls d1*

- ▶ *ls d1 / d1*

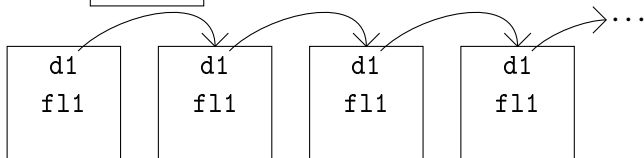
- ▶ *ls d1 / d1 / d1/*

Smart sharing (cont)

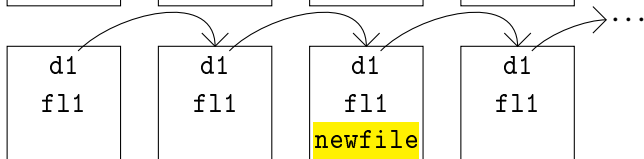
Initially, our file system was like that



However, it appeared like this



After we created `/d1/d1/newfile`, it became:



When we updated the directory `/d1 / d1`, the zipper automatically broke the cycle and introduced three real directories. We get the real copy-on-write. Try to get this with a Unix file system!

Conclusions

Zipper-based file system over any term

- ▶ Transactional semantics
- ▶ Strongest (repeatable read) isolation mode
- ▶ Built-in traversal
- ▶ Smart sharing
- ▶ Threading and exceptions via delimited continuations
- ▶ Static guarantee of processes' non-interference

Future work

- ▶ FUSE or NFS or 9P server
- ▶ Semantically richer terms: extended attributes, ...
- ▶ *cd* into a λ -term in *bash*

Delimited continuations do the right thing for free