

# Forest 1.0: A Language and Toolkit For Programming with Filestores

Kathleen Fisher  
Tufts University

Nate Foster  
Cornell University

David Walker  
Princeton University

Kenny Q. Zhu  
Shanghai Jiao Tong University

## Abstract

A *filestore* is a structured collection of data files housed in a conventional hierarchical file system. Many applications use filestores as a poor-man’s database, and the correct execution of these applications requires that the collection of files, directories, and symbolic links stored on disk satisfy a variety of precise invariants. Moreover, all of these structures must have acceptable ownership, permission, and timestamp attributes. Unfortunately, current programming languages do not provide support for documenting assumptions about filestores, detecting errors, or safely loading from and storing to them.

This paper describes the design, implementation, and semantics of Forest, a novel domain-specific language for describing filestores. The language uses a type-based metaphor to specify the expected structure, attributes, and invariants of filestores. Forest generates loading and storing functions that make it easy to connect data on disk to an isomorphic representation in memory that can be manipulated as if it were any other data structure. Forest also generates metadata that describes the degree to which the structures on the disk conform to the specification, making error detection easy. Hence, in a nutshell, Forest extends the rigorous discipline of typed programming languages and many of their benefits to the untyped world of file systems.

We have implemented Forest as an embedded domain-specific language in Haskell. In addition to generating infrastructure for reading, writing and checking file systems, our implementation generates a type class instances that make it easy to build generic tools that operate over arbitrary filestores. We illustrate the utility of this infrastructure by building a file system visualizer, a file access checker, a generic query interface, description-directed variants of several standard UNIX shell tools and (circularly) a simple Forest description inference engine. Finally, we formalize a core fragment of Forest in a semantics inspired by classical tree logics and prove round-tripping laws showing that the loading and storing functions behave sensibly.

## 1. Introduction

Databases are an effective, time-tested technology for storing structured and semi-structured data. Nevertheless, many computer users eschew the benefits of structured databases and store important semi-structured information in collections of files and directories in a conventional file system instead. For example, the Princeton Computer Science Department stores records of undergraduate student grades in a structured set of directories and uses scripts to compute averages and study grading trends. Similarly, Michael Freedman collects sets of log files from CoralCDN, a distributed content distribution network [11, 12]. The logs are organized in hierarchical directory structures based on machine name, time and date. Freedman mines the logs for information on system security and performance. At Harvard, physics professor Vinodhan Manoharan stores his experimental data in sets of files and extracts information using python scripts. At AT&T, vast structured repositories contain net-

work monitoring information, phone call records, and billing data. Many software code bases, including Haskell and its associated Cabal libraries, require that specific files exist in particular formats at precise locations described in other files. Similarly, version control systems like `cvs` utilize the file system to store revision information. Web sites require various types of files to exist in particular directories according to content type, and security considerations often require particular permissions on these files. Many other examples exist across the computational sciences and social sciences, in computer systems research, in computer systems administration and in industry.

Users choose to implement ad hoc databases in this manner for a number of reasons. A key factor is that using databases often requires paying substantial up-front costs such as: (1) finding and evaluating the appropriate database software (and possibly paying for it); (2) learning how to load data into the database; (3) writing programs to transform the raw data for loading into the database; (4) learning how to access the data once it is in the database; and (5) interfacing the database with a conventional programming language to support applications that use the data. Finally, it may be the case that the database optimizes for a pattern of use not suited to the actual application, which makes the overhead of the database system even less desirable.

Rather than paying these costs, programmers often store data in the file system, using a combination of directory structure, file names, file contents, and symbolic links to organize the data. We call such a representation of a coherent set of data a *filestore*. The “query language” for a filestore is often a shell script or conventional programming language.

Unfortunately, despite their initial convenience, using filestores can have a number of negative consequences. First, there is generally no documentation, which means it can be hard to understand the data and its organization. New users struggle to learn the structure, and if the system administrator leaves, knowledge of the data organization may be lost. Second, the structure of the filestore tends to evolve: new elements are added and old formats are changed, sometimes accidentally. Such evolution can cause hacked-up data processing tools to break or return erroneous results; it also further complicates understanding the data. Third, there is often no systematic means for detecting errors even though data errors can be immensely important. For example, for filestores containing monitoring information, errors can signal that some portion of the monitored system is broken. Fourth, analyses tend to be built from scratch. There is no auxiliary query or tool support and no help with debugging. Tools tend to be “one-off” efforts that are not reuseable. Fifth, dealing with large data sets, which are common in this setting, imposes extra difficulties. For example, standard shell tools such as `ls` fail when more than 256 files appear on the command line. Hence, programmers must break up their data and process it in smaller sets, a tedious task.

We propose a better way: A type-based specification language, programming environment and toolkit for describing and managing filestores. This language, called Forest, is implemented as an

embedded domain-specific language in Haskell. Forest allows programmers to describe the expected shape of a filestore and to materialize it as typed, format-specific Haskell data structures. Conversely, given data structures with the appropriate type, Forest makes it straightforward to dematerialize these structures and write them out to disk.

The first benefit of the Forest system is that Forest descriptions provide *executable documentation* that can be used to check whether a given filestore conforms to its specification. For example, Unix file systems should be laid out according to the informal standard set forth by the Filesystem Hierarchy Standard Group [3], which requires, among other things, that certain directories *only* contain certain files, presumably for security reasons. Forest provides a language for expressing standards precisely and for checking that given file systems conform to them. As another example, the Pads website [25] contains a complex set of scripts and data files to implement an online demo. Unless all of the required data files, directories, and symbolic links are configured correctly, the web demo fails with an inscrutable error message. Forest allows the Pads webmaster to precisely document all of these requirements and to detect specification violations, making it easy to find and repair errors. And, of course, if the current webmaster were to leave her post, her successor could use the Forest description to help understand the system.

As well as serving as executable documentation, Forest provides substantial additional support for programmers. The goal is for programmers to obtain a whole range benefits by writing one simple, compact file system specification. The automatically generated auxiliary support includes: (1) a set of type declarations to represent the filestore in memory; (2) a set of type declarations that capture errors and file system attributes for the filestore; (3) a loading function to populate these in-memory structures; (4) a storing function to push possibly updated structures back out to disk; (5) type class instance declarations that make it possible for programmers to query, analyze, and transform filestore data using generic functions; and (6) a set of useful generic functions/scripts that operate over instances of these type classes.

Overall, the main contribution of this work is conceptual: We propose the idea of extending a modern programming language with tightly integrated linguistic features for describing and manipulating filestores. To demonstrate the potential of this idea, the following sections of this paper flesh our proposal in greater depth:

- Section 2 begins with two concrete motivating examples, drawn from the authors day-to-day experience managing computer systems. While there are just two central examples in this paper, the Forest web site [9] contains a number of further examples and case studies.
- Section 3 describes a concrete language design. The design is characterized by a simple, intuitive and compositional syntax that is tightly integrated with Haskell, our host language. The design is also tightly integrated with Pads/Haskell, a domain-specific language for describing individual files (as opposed to entire file stores), inspired by past work on related data description languages [5, 6, 7, 22]. This tight, seamless integration was a crucial design goal as it allows programmers to transition effortlessly between ordinary Haskell data structures, file internals and file collections, all in a uniform syntax.
- Section 4 explains how to write Haskell programs that operate over filestores described in Forest. The goal of this section is to provide a sense of just how easy it is to write simple file system scripts or queries.
- Section 5 shows that it is possible to use Forest to make management of filestores even easier by developing general-

purpose, generic tools capable of operating over any filestore! We have developed several of these tools including a generic query interface, a file system visualization tool, an access control permission checker, and a series of UNIX-like scripting tools. We have also built a simple description-inference tool to help users write a new description for a given an existing file system. These tools are interesting in their own right and also as case studies of putting generic programming techniques into practice. In addition, they provide evidence that our design is effectively integrated into the Haskell ecosystem.

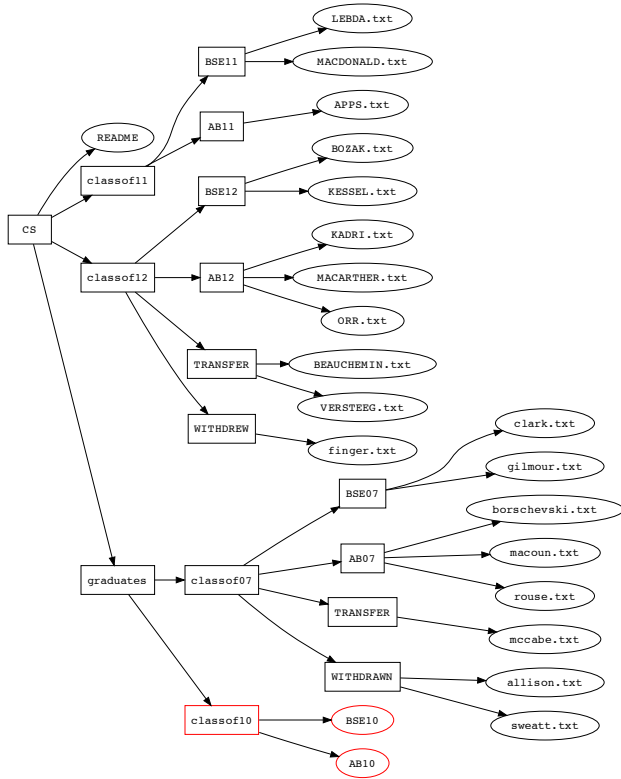
- Section 6 explains our implementation, which is complete and may be downloaded at the Forest web site [9]. In addition to delivering a useful tool, our engineering work here has the auxiliary benefit of serving as a rich case study in domain-specific language implementation. In fact, it has already had significant impact as such by influencing the development of Haskell itself: the Haskell team modified and extended the quasiquoting mechanism in response to our needs.
- Section 7 describes the formal semantics for core Forest and states theorems demonstrating that the mappings between the file system and in-memory structures behave correctly. These theorems are inspired by the "round-tripping" laws for well-behaved lenses [10], but are significantly more complicated as the load and store functions have to deal with inconsistencies stemming from dependencies, duplication, and invalid data.
- Section 8 contains a discussion of related work. There has been much past work on domain-specific languages for describing, parsing and printing individual data files. Examples include Lex, Yacc, Antlr [26], Parsec [19] and Pads [7], to name just a few. However, Forest differs substantially from any of these systems because it focuses on technology for describing *entire file systems*. A key difference is that simple file systems are *trees* and complex ones with symbolic links are *graphs*, whereas files are *sequences* (of characters or tokens). Consequently, the language design, formal systems, semantic issues, and underlying implementation technology are all entirely different.
- Finally, Section 9 concludes.

## 2. Example Filestores

In this section, we present two example filestores. We use these examples to motivate and explain the design of Forest.

The first filestore contains information about students in Princeton's undergraduate computer science program. The faculty use the information to decide on undergraduate awards and to track grading trends. Its format has changed over time—something that is typical for ad hoc filestores! Naturally, any description needs to cope with the variations introduced as formats evolve.

Figure 1 shows a snippet of the (anonymized) student filestore designed to illustrate its structure. At the top level, there are three directories: `classof11` (seniors), `classof12` (juniors) and `graduates` (students who have graduated). There is also a `README` file containing a collection of notes. Inside `graduates`, there is set of directories named `classofYY` where `YY` dates back to 92. Inside each `classofYY` directory, there are at least the two degree subdirectories `ABYY` and `BSEYY` as the computer science department gives out both Arts and Science (AB) and Engineering (BSE) degrees. Optionally, there are also subdirectories for students who withdrew from Princeton or transferred to another program. Within any degree subdirectory, there is one text file per student that records the courses taken and the corresponding grades. Each degree directory may also contain a template file named `sss.txt` or `sxx.txt` for creating new students.



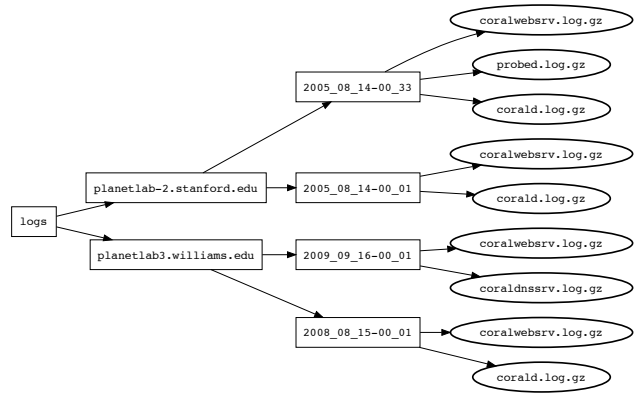
**Figure 1.** Anonymized snippet of Princeton computer science undergraduate data. Red notes denotes an error—i.e., missing files.

The second filestore contains log files for CoralCDN [11, 12]. To monitor the performance and security of the system, the hosts participating in CoralCDN periodically send usage statistics back to a central server. These statistics are collected in a filestore similar to the one depicted in Figure 2. The filestore has a top-level directory named `dat`, which contains a set of subdirectories, one for each host. Each of those directories contain another set of directories, labeled by date and time. Finally, each of the date/time directories contain one or more compressed log files. For the purposes of this example, we will focus on the `coralwebserv.log.gz` log file, which contains detailed information about the web requests made on the host during the preceding time period. In addition to exploring this primary filestore, we also explore a secondary, derived filestore. This secondary store, named `stats`, contains files that store statistics generated by Forest/Haskell scripts that analyze and summarize the raw CoralCDN server data. These system-wide summaries are representative of the statistics reported by Freedman in his CoralCDN report [11].

### 3. Forest Design

Data stored in a filestores shares many characteristics of data stored in ordinary, in-memory programmer data structures. Consequently, Forest uses the same sort of language to describe filestores as one uses to describe ordinary data structures — the language of types. Simple base types describe individual file system objects<sup>1</sup> and more complex types describe organized collections of file system objects. This idea forms the basis for our design.

<sup>1</sup>We use the term “file system object” or more simply “object” to denote either a file, or a directory, or a symbolic link.



**Figure 2.** Coral system log data.

**Embedding Forest in Haskell.** In order to write lightweight scripts, programmers must be able to manipulate and transform file system objects side-by-side with ordinary data structures. Consequently, a language like Forest must be embedded within a more general host programming language. We chose Haskell as the host language primarily because of its rich support for type-directed programming, which facilitates the construction of generic tools that can operate over any Forest description. As a bonus, Haskell’s quasiquoting mechanism [21] proved a useful way to implement Forest. It enabled tight integration of the two languages, while admitting fine-grained control over Forest syntax.

To introduce new Forest declarations within a Haskell program, the programmer simply opens the Forest sublanguage using quasiquoting notation:

```
[forest| ... forest declarations ... |]
```

When processing such a quasiquote, the Haskell compiler invokes the Forest compiler, which converts the given Forest declarations into a sequence of plain Haskell declarations that collectively implement the Forest declarations.

**Forest Structure and Interpretations.** Once within the Forest sublanguage, the programmer writes declarations that resemble extended Haskell type declarations. Each such type declaration has three primary semantic interpretations:

1. An interpretation as an *expected on-disk shape* of a file system fragment.
2. An interpretation as an ordinary Haskell type for the in-memory *representation* that will be constructed when the file system fragment is loaded into a Haskell program.
3. An interpretation as an ordinary Haskell type for the in-memory *metadata* that will be generated when the file system fragment loaded.

All three interpretations are used by the tool that loads data from the file system into memory as specified by a Forest description. When supplied with a *current path*, the loader uses the first interpretation to validate that the file system rooted at that path has the correct shape. If the expected shape is complicated, possibly involving several nested subshapes (and hence traversal through several subdirectories), the semantics of Forest dictates how the loader should adjust the current path as it goes. When validation (also called *matching*) succeeds, we say the file system fragment *matches* the description. The second interpretation is used when the loader lazily pulls the on-disk data into memory. The in-memory data structure is guaranteed to have the Haskell type given by the second interpretation. The third interpretation provides a type for the meta-

data structure generated by the loader. Such metadata includes error information (missing file, insufficient permissions, *etc.*) as well as file system attributes (owner, size, *etc.*).

The effectiveness of the Forest language comes in part from the fact that these three interpretations all arise from a single compact description. Moreover to aid the programmer in navigating between interpretations, we align the syntax of Forest with the syntax of Haskell where possible. For example, if the Haskell types for the in-memory representation and metadata are record types then the Forest syntax is designed to look similar to a Haskell record type. Likewise, if the Haskell types for the in-memory representation and metadata are `Maybe` types then the Forest syntax is designed to look similar to a Haskell `Maybe` type. Many of these high-level design considerations were adopted from earlier work on Pads [5, 7, 22], though, as mentioned earlier, the semantics of Forest (which operates over graph-based file systems) is substantially different from the semantics of Pads (which operates over sequence-based strings).

**Errors.** As with Pads [5, 7, 22], we do not assume that a given filestore conforms perfectly to its associated Forest description. Instead, we check during the loading process the extent to which the filestore conforms, marking discrepancies in the metadata. This design allows users to respond in application-specific ways to errors. It also allows us to check dynamically the conditions implied by Forest’s dependent types, skirting issues of undecidability.<sup>2</sup> Because Forest loads data lazily, this choice means errors will not be detected unless the user program needs to touch the portion of the filestore with the error. The user can force a complete conformance check by accessing the top-level error count. It is possible for the filestore to change during or after this check. For the filestores we have seen in practice, there are extra-linguistic procedures in place to prevent such concurrent modifications; we leave to future work the possibility of using operating system support to monitor and/or prevent such changes automatically.

**Onward.** In the remainder of this section, we discuss the specific type constructors that constitute the Forest language and illustrate their use in our running examples.

### 3.1 Base Types: Files

Forest provides a small collection of base types for describing individual files: `Text` for ASCII files, `Binary` for binary files, and `Any` for arbitrary files. As with all Forest types, each of these types specifies a representation type, a metadata type, and loading and storing functions. For all three file types, the representation type is a `ByteString`. Similarly all three share a metadata type, which pairs file-system metadata with metadata describing properties of the file contents. The file-system metadata has type `Forest_md`, shown in Figure 3. This structure stores two kinds of information:

1. the number and kind of any errors that occurred during loading
2. the attributes associated with the file (`fileInfo`)

File-content metadata describes errors within the file. For these three file types, there is no meaningful content metadata and so this type is the unit type. Leveraging Haskell’s laziness, the loading functions create the in-memory representations and set the metadata on demand. The storing functions, which are described in more detail in Section 4, do the inverse.

<sup>2</sup>Validation that a file system obeys a Forest specification is akin type checking. However, it is akin to type checking first-order values (trees and graphs) as opposed to type checking higher-order values (functions). Consequently, even though Forest has dependent types, type checking is a simple linear traversal of the file system. Forest does not have to decide equivalence of expressions with free variables as one must do when type checking a dependent lambda calculus, for example.

```

data Forest_md = Forest_md
  { numErrors :: Int
  , errorMsg  :: Maybe ErrMsg
  , fileInfo  :: FileInfo }

data FileInfo = FileInfo
  { fullpath  :: FilePath
  , owner     :: String
  , group     :: String
  , size      :: Coff
  , access_time :: EpochTime
  , mod_time  :: EpochTime
  , read_time :: EpochTime
  , mode      :: FileMode
  , isSymLink :: Bool
  , kind      :: FileType }

```

Figure 3. Forest metadata types.

Although useful, these three base types are not sufficient for describing the wide range of files used in practice, including XML documents, Makefiles, source files in various languages, shell scripts, *etc.* The appropriate representation and content metadata types for each such file varies. To support such files, Forest provides a plug-in architecture, allowing third-party users to define new file types by specifying a representation type, a metadata type, and corresponding loading and storing functions.

A common class of files are *ad hoc data files* containing semi-structured information, an example of which is the Princeton student record file format. In such cases, Forest can leverage the Pads/Haskell [8] data description language to define format-specific in-memory representations, content metadata, and loading and storing functions. Pads/Haskell is a recently developed version of Pads [5, 7, 22]. Like Forest, Pads/Haskell is embedded in Haskell using quasiquotation. For example, the following code snippet begins the Pads specification of the Princeton student record format:

```

[pads| data Student (name :: String) =
  { person :: Line (Person name)
  , Header
  , courses :: [Line Course]
  , Trailer
  }
... |]

```

This description is parameterized by the name of the student whose data is in the file; the complete description appears in the companion technical report [4]. From this specification, the Pads compiler generates an in-memory representation type `Student`, a content metadata type `Student_md`, and parsing and printing functions.

Forest provides the `File` type constructor to lift Pads types to Forest file types. For example, the declaration

```

[forest| type SFile (n :: String) = File (Student n) |]

```

introduces a new file type named `SFile` whose format is given by the Pads type `Student`. As with the Pads type, `SFile` is parameterized by the name of the student.

Using Pads/Haskell descriptions in Forest not only helps specify the structure of ad hoc data files, but it also generates a structured in-memory representation of the data, allowing Haskell programmers to traverse, query and otherwise manipulate such data. We designed Pads/Haskell and Forest to work seamlessly together. From the perspective of the Haskell programmer traversing a resulting in-memory data structure, there is effectively no difference between iterating over files in a directory or structured sequences of lines or tokens within a file.

While Pads/Haskell is independently interesting, this paper focuses on Forest. Henceforth, any unadorned declarations occur

within the Forest scope `[forest | . . . |]` unless otherwise noted. Any declarations prefixed by `>` are ordinary Haskell declarations.

### 3.2 Base Type: Symbolic Links

When symbolic links occur in a described filestore, Forest follows the symbolic link to its target, mimicking standard shell behavior. However, Forest allows programmers to specify explicitly that a particular file is a symbolic link using the base type `SymLink`. The in-memory representation for an explicit symbolic link is the path that is the target of the link. It is possible to use constraints (Section 3.6) to specify desired properties of the link target, such as requiring it to be to a specific file.

In Forest, any file system object may be described in multiple ways. Hence, in the case of a symbolic link, it is possible to use one declaration to specify that the object is a symbolic link and a second to specify the type of the link target. We will see an example of such a specification in the next subsection.

### 3.3 Maybe: Optional File System Objects

Sometimes, a given file (or directory or symbolic link) may or may not be present in the file system, and either case is valid. To handle this situation, we leverage the idea of an option type by providing a Forest-level `Maybe` type constructor that maps the optional file system object to a `Maybe` type in Haskell. In particular, if `T` is a Forest type, then `Maybe T` is the Forest type denoting an optional `T`. The type `Maybe T` succeeds and returns representation `None` when the current path does not exist in the file system. `Maybe T` also succeeds and returns `Just v` for some `v` of type `T` when the current path exists and matches `T`. `Maybe T` registers an error in the metadata when the current path exists but the corresponding object does not match `T`.

### 3.4 Records: Directories

Forest directories are record-like datatype constructors that allow users to specify directory structures. For example, to specify the root directory of the student repository in Figure 1, we might use the following declaration. This declaration assumes that we have already defined `Class y`, a parameterized description that specifies the structure of a directory holding data for the class of year `y`, and `Grads`, a description that specifies the structure of the directory holding all graduated classes.

```
type PrincetonCS_1 = Directory
{ notes is "README" :: Text
, seniors is "classof11" :: Class 11
, juniors is "classof12" :: Class 12
, grads is "graduates" :: Grads }
```

Each field of the record describes a single file system object. It has three components: (1) an internal name (e.g., `notes` or `seniors`) that must be a valid Haskell record label, (2) an external name specified as a value of type `String` (e.g., `"README"` or `"classof11"`) that gives the name of the object on disk, and (3) a Forest description of the object (e.g., `Text` or `Class 11`).

When the external name is itself a valid Haskell label, users may omit it, in which case Forest uses the label as the on-disk name:

```
type PrincetonCS_2 = Directory
{ notes is "README" :: Text
, classof11 :: Class 11
, classof12 :: Class 12
, graduates :: Grads }
```

We could not abbreviate the `notes` field because labels must start with a lowercase letter in Haskell.

**Matching.** For a file system object to match a directory description, the object must be a directory and each field of the record must

match. A field `f` matches when the object whose path is the concatenation of the current path and the external name of `f` matches the type of `f`.

It is possible for the same file system object to match multiple fields in a directory description at the same time. For example, if `"README"` were actually a symbolic link, it is possible to document that fact by mentioning it twice in the directory description, once as a text file and once as a symbolic link:

```
type PrincetonCS_3 = Directory
{ link is "README" :: SymLink
, notes is "README" :: Text
, ... }
```

It is also possible for a directory to contain objects that are unmatched by a description. We allow extra items because it is common for directories to contain objects that users do not care about. For example, a directory structure may contain extra files or directories related to a version control system, and a description writer may not want to clutter the Forest specification with that information. We will see shortly that it is possible to specify the absence of file system objects using constraints.

As suggested by the syntax, the in-memory representation of a directory is a Haskell record with the corresponding labels. The type of each field is the representation type of the Forest type for the field. The metadata has a similar structure. The metadata for each field has two components: file-system attribute information of type `Forest_md` and field-specific metadata whose type is derived from the Forest type for the field. In addition, the directory metadata contains an additional value of type `Forest_md` that summarizes the errors occurring in directory components and stores the `FileInfo` structure for the directory itself. When loading a directory, Forest constructs the appropriate in-memory representation for each field that matches and puts the corresponding metadata in the metadata structure. For fields that do not match, Forest constructs default values and marks the metadata with suitable error information.

**Computed Paths** The above descriptions are a good start for our application, but they are not ideal. Every year, the directory for graduating seniors (i.e., `classof11`) is moved into the graduates directory, the juniors are promoted to seniors and a new junior class is created. As it stands, we would have to edit the description every year. An alternative is to parameterize the description with the current year and to *construct* the appropriate file names using Haskell functions:

```
> toStrN i n = (replicate(n - length(show i)) '0')
> ++ (show i)
> mkClass y = "classof" ++ (toStrN y 2)
```

```
type PrincetonCS (y::Integer) = Directory
{ notes is "README" :: Text
, seniors is <|mkClass y |> :: Class y
, juniors is <|mkclass (y+1)|> :: Class <|y+1|>
, graduates :: Grads }
```

The bracket syntax `<| . . . |>` provides an escape so that we may use Haskell within Forest code to specify arbitrary computations. When an expression is a constant or variable, it may be supplied directly. When an argument is more complex, however, it must be written in brackets to escape to Haskell. This example also illustrates abstraction: any Forest declaration may be parameterized by specifying a legal Haskell pattern and its type. The types of the fields for `seniors` and `juniors` illustrate the use of parameterized descriptions.

**Approximate Paths** As filestores evolve, naming conventions may change. Additionally, directory structures with multiple instances may have minor variations in the names of individual files

across instances. For example, in each Princeton class directory, there may (or may not) be some number of students that have withdrawn from the program, transferred to a different program, or gone on leave. Over the years, slightly different directory names have been used to represent these situations.

To accommodate this variation, Forest includes the matching construct to approximate file names. We can use this mechanism to describe the class directory:

```
> transRE = RE "TRANSFER|Transfer"
> leaveRE = RE "LEAVE|Leave"
> wdRE     = RE "WITHDRAWN|WITHDRAWAL|Withdrawn"

type Class (y::Integer) = Directory
  { bse is <|"BSE" ++ (toString y)|> :: Major
  , ab  is <|"AB"  ++ (toString y)|> :: Major
  , trans matches transRE :: Maybe Major
  , withd matches wdRE    :: Maybe Major
  , leave matches leaveRE :: Maybe Major }
```

A field with the form `<label> matches <regex> :: T` finds the set of paths in the files system that match `currentPath/<regex>`. If there are zero or one such files, the `matches` form acts just as the `is` form. If more than one file matches, one of the matches is selected non-deterministically, a multiple match error is registered in the metadata, and matching continues as it would with the `is` form. In addition to regular expressions, the matching construct also allows *glob patterns*, (i.e., patterns such as `*.txt`), to specify the names of files on disk. An example appears in the next subsection.

### 3.5 Lists

Just as Haskell has both records and lists, so too does Forest. Records allow programmers to specify a fixed number of file system objects, each with its own type. Lists, on the other hand, allow programmers to specify an arbitrary number of file system objects, each with the same type. As an example, we can use a list to specify the Grads directory from Figure 1. We borrow Haskell's notation for list comprehensions to specify the names of the file system objects:

```
> getYear s =
>   toInteger $ reverse $ take 2 $ reverse s
> cRE = RE "classof[0-9][0-9]"

type Grads =
  [c :: Class <|getYear c> | c <- matches cRE]
```

In this specification, `Grads` is a directory fragment containing a number of `Class` subdirectories with names `c` that match the regular expression `cRE`. The Haskell function `getYear` extracts the last two digits from the name of the directory, converts the string digits to an integer year, and passes the year to the underlying `Class` specification. More generally, Forest lists have the form `[path :: T | id <- gen, pred]` where `id` is bound in turn to each of the file names generated by `gen`, which may be a `matches` clause (used to match against the files at the current path as in the previous section) or a list computed in Haskell. These generated `ids` are filtered by the optional predicate `pred`. For each such legal `id`, there is a corresponding expression `path`, which Forest interprets as extending the current path. The object at each such path should have the Forest type `T`. The identifier `id` is in scope in `pred`, `path`, and `T`.

The in-memory representation of a Forest list is a Haskell list containing pairs of the name of a matching object and its representation. The metadata is a list of the metadata of the matching objects paired with a summary metadata structure of type `Forest_md`.

**Representation Transformations.** Although the list representation for comprehensions is useful, it can be desirable to use a more sophisticated data structure to represent such collections. To support this usage, Forest allows programmers to prefix a list comprehension with any type constructor that belongs to a Forest-defined container type class. This type class contains functions that specify how to convert between the list representation and the desired container representation. We have provided such instance declarations for Haskell's `Map` and `Set` type constructors.

As an example, consider the specification of the `Major` directory. Each such directory contains a list of student files and an additional template file named either `sss.txt` or `sxx.txt`. The declaration below specifies the collection of student files by matching with a glob pattern and filtering to exclude template files. It uses the `Map` type constructor to specify that the data and metadata should be collected in a `Map` rather than a list.

```
> template s = s `elem` ["sss.txt", "sxx.txt"]
> txt = GL "*.txt"
```

```
type Major = Map
  [ s :: File (Student <|dropExtension s|>)
  | s <- matches txt, <|not (template s)|>]
```

### 3.6 Dependent Types: Attributes and Constraints

Every file system object has a number of *attributes* associated with it, such as its owner, group, permissions, and size. In general, if a Forest identifier `id` refers to a path, then the identifier `id_att` refers to the corresponding attributes. This attribute identifier has type `Forest_md`, shown in Figure 3. Forest defines helper functions to access these attributes, some of which are listed in Figure 4.

*Constrained types* are a simple form of dependent types that allow users to specify required attributes. For example, the type `PrivateFile` specifies a text file accessible only by its owner.

```
type PrivateFile =
  Text where <|get_modes this_att == "-rw-----"|>
```

The keyword `where` introduces a constraint on the underlying type. The load function for the type `PrivateFile` checks this constraint during loading. If the constraint is false, it records that fact in the metadata. Within constraints, the special identifier `this` refers to the representation of the underlying object, `this_att` refers to its attributes and `this_md` to its complete metadata.

Using attributes, we can write a *universal directory description*, which is sufficiently general to describe any directory:

```
type Universal = Directory
  { asc is [ f :: Text
            | f <- matches (GL "*"),
            <| get_kind f_att == AsciiK |> ]
  , bin is [ b :: Binary
            | b <- matches (GL "*"),
            <| get_kind b_att == BinaryK |> ]
  , dir is [ d :: Universal
            | d <- matches (GL "*"),
            <| get_kind d_att == DirectoryK |> ]
  , sym is [ s :: SymLink
            | s <- matches (GL "*"),
            <| get_sym s_att == True |> ] }
```

When a directory is loaded using the `Universal` description, all the ASCII files will end up the `asc` field, all the binary files in `bin`, all the directories in `dir`, and all the symbolic links in `sym`. Note that the description uses recursion to describe nested directories. In the case that a symbolic link creates a cycle in the file system by pointing to a parent directory, the Haskell in-memory representation is a (lazy) infinite data structure. We view the fact that it is possible to write such a universal description in

function name	information
get_group	object group
get_kind	the sort of file or directory
get_modes	permission string
get_owner	object owner
get_size	object size

Figure 4. Selected file attribute functions

```
[forest|
data PrincetonCS (y::Integer) = Directory
  { notes is "README" :: Text
  , seniors is <|mkClass y |> :: Class y
  , juniors is <|mkClass (y+1)|> :: Class <|y+1|>
  , graduates :: Grads }

data Class (y::Integer) = Directory
  { bse is <|"BSE" ++ (toString y)|> :: Major
  , ab is <|"AB" ++ (toString y)|> :: Major
  , trans matches transRE :: Maybe Major
  , withd matches wdRE :: Maybe Major
  , leave matches leaveRE :: Maybe Major }

type Grads =
  [ c :: Class <|getYear c|> | c <- matches cRE ]

type Major = Map
  [ s :: File (Student <|dropExtension s|>)
  | s <- matches txt, <|not (template s)|> ] ]
```

Figure 5. Forest description of Princeton filestore.

Forest as evidence that the language is appropriately expressive. This description also serves as an example of how to describe a filestore by its *structure* rather than its *names*.

We can also use constraints to specify that certain files *do not* appear in certain places. As an example, we might want to require that no binaries appear in a directory given to an untrusted user as scratch space. The description below flags an error during loading if a binary file exists in the directory.

```
type NoBin =
  [ b :: Binary | b <- matches (GL "*"),
    <| get_kind b_att == BinaryK |> ]
where <|length this == 0|>
```

### 3.7 Specialized Constructors: Gzip and Tar

Some files need to be processed before they can be used. A typical example is a compressed file such as the gzipped log files in CoralCDN. Forest provides processing-specific type constructors to describe such files. For example, if `CoralLog` is a Pads/Haskell description of a CoralCDN log file then

```
type Info = Gzip (File CoralLog)
```

describes a gzipped log file. Likewise, suppose `logs.tar.gz` is a gzipped tar file and that the type `ManyLogs` describes the directory of log files that `logs.tar` expands to when untarred. Such a situation can be described using a combination of the `Tar` and `Gzip` type constructors:

```
type MoreInfo = Gzip (Tar ManyLogs)
```

### 3.8 Putting it all together

The preceding subsections give an overview of Forest's design. Figures 5 and 6 give the specifications for the two running examples,

```
[forest|
type Stats = Directory
  { last :: File Last, topk :: File Topk }
type Dat = [ s :: Site | s <- matches site ]
type Site = [ d :: Log | d <- matches time ]
data Log = Directory
  { log is coralwebsrv :: Gzip (File CoralLog) } ] ]
```

Figure 6. Forest CoralCDN description.

### Coral Representation Types:

```
newtype Stats = Stats {last :: Last, topk :: Topk}
newtype Dat = Dat [(String, Site)]
newtype Site = Site [(String, Log)]
data Log = Log {log :: CoralLog}
```

### Coral Metadata Types:

```
type Stats_md = (Forest_md, Stats_inner_md)
data Stats_inner_md = Stats_inner_md
  {last_md :: (Forest_md, Last_md),
  topk_md :: (Forest_md, Topk_md)}
type Dat_md = (Forest_md, [(String, Site_md)])
type Site_md = (Forest_md, [(String, Log_md)])
type Log_md = (Forest_md, Log_inner_md)
data Log_inner_md = Log_inner_md
  {log_md :: (Forest_md, CoralLog_md)}
```

### Load Functions:

```
stats_load :: FilePath -> IO (Stats, Stats_md)
dat_load :: FilePath -> IO (Dat, Dat_md)
site_load :: FilePath -> IO (Site, Site_md)
log_load :: FilePath -> IO (Log, Log_md)
```

### Store Functions:

```
stats_manifest :: (Stats, Stats_md) -> IO Manifest
dat_manifest :: (Dat, Dat_md) -> IO Manifest
site_manifest :: (Site, Site_md) -> IO Manifest
log_manifest :: (Log, Log_md) -> IO Manifest
storeAt :: FilePath -> Manifest -> IO ()
store :: Manifest -> IO ()
```

Figure 7. Coral rep. and metadata types; load and store functions

minus the associated Pads/Haskell and Haskell declarations. The complete descriptions of these filestores and additional descriptions are available in a technical report [4], including descriptions of the Pads website, a Gene Ontology filestore, and CVS repositories.

## 4. Programming with Forest

Many Forest programs work in two phases. In the first phase they use Forest to load relevant portions of the file system into memory, and in the second phase they use an ordinary Haskell function to traverse the in-memory representation of the data (or its associated metadata) and compute the desired result. Some Forest programs add a third phase in which they store updated structures back to the filestore.

To facilitate this style of programming, the Forest compiler generates several Haskell types and functions from every Forest declaration. Collectively, these types and functions define an instance of the `Forest` type class:

```
class (Data rep, ForestMD md)
=> Forest rep md | rep -> md where
  load :: FilePath -> IO(rep, md)
  manifest :: (rep,md) -> IO Manifest
  ...
```

In this type class, the type `rep` is the generated in-memory representation type of the corresponding on-disk data. The type `md` is the generated type for the associated metadata. The `ForestMD` type class provides operations for manipulating Forest metadata; all generated metadata types belong to this type class.

The generated load function lazily traverses the file system and reads the files, directories, and symbolic links mentioned in the description into a pair of the in-memory representation and its metadata. To reverse the process of reading data in to memory, Forest also generates a *manifest function*, which reads an in-memory data structure, writes its contents out to disk in a temporary space, and prepares a *manifest log*. The manifest log records inconsistencies detected during this process as well as the sequence of operations necessary to move data from the temporary space to its final resting point. Inconsistencies can arise when a programmer creates an erroneous in-memory representation of a filestore. The dependencies that may be present in Forest descriptions mean that not all such inconsistencies can be detected statically by the Haskell type system. After creating a manifest, a programmer may analyze it and decide whether to execute the generic `store` or `storeAt` functions, which move a manifest (inconsistencies and all) to its rightful position on disk. Details concerning the semantics of storing, especially where it concerns inconsistencies, are explained in further depth in Section 7.

As an example, consider the CoralCDN logs described in Figure 6. The corresponding load and store functions, the representation types, and the metadata types appear in Figure 7.<sup>3</sup> Note that the structure of each of these artifacts mirrors the structure of the Forest description that generated them. This close correspondence makes it easy for programmers to write programs using these Forest-generated artifacts.

For instance, consider the `Dat` description in Figure 6. The `dat_load` function takes a path as an argument and produces the representation and metadata obtained by loading each of the site directories contained in the directory at that path:

```
(rep,md) <- dat_load "/var/log/coral/dat"
```

Because `Dat` is a Forest list, the `rep` is a Haskell list. More specifically, `rep` has the form

```
Coral [ ("planetab2.eecs.wsu.edu", Site [...]),
       ("planetlab3.williams.edu", Site [...]), ... ]
```

where the list contains pairs of names of subdirectories and representations for the data loaded from those directories. The metadata is a pair consisting of a generic header of type `Forest_md` and a list of pairs of names of subdirectories and their associated metadata. The header collects information about errors encountered during loading and it stores the file system attributes of each file, directory, or symbolic link loaded from the file system. The following is the pretty-printed version of such a structure:

```
Forest_md
{ numErrors = 0,
  errorMsg = Nothing,
  fileInfo = FileInfo
  { fullpath = /var/log/coral/dat,
    owner = alice, group = staff, size = 102,
    access_time = Fri Nov 19 01:47:09 2010,
    mod_time = Thu Nov 18 20:42:37 2010,
    read_time = Fri Nov 19 01:47:28 2010,
    mode = drwxr-xr-x, isSymLink = False,
    kind = Directory } },
[ ("planetlab2.eecs.wsu.edu", Forest_md {...}),
  ("planetlab3.williams.edu", Forest_md {...}), ... ]
```

<sup>3</sup>In the following examples, for the sake of clarity, we use type-specific names such as `dat_load` and `dat_manifest`, rather than the overloaded names `load` and `manifest`.

Using these functions and types, it is easy to formulate many useful queries as simple Haskell programs. For instance, to count the number of sites we can simply compute the length of the nested list in `rep`:

```
num_sites = case rep of Dat l -> List.length l
```

More interestingly, since the internals of the web log are specified using `Pads/Haskell` (see the technical report [4] for details), it is straightforward to dig in to the file data and combine it with file metadata or attributes in queries. For example, to calculate the time when statistics were last reported for each site, we can zip the lists in `rep` and `md` together and project out the site name and the `mod_time` field from each element in the resulting list of pairs:

```
get_site = fst
get_mod (_, (f,_)) = mod_time . fileInfo $ f
sites_mod () =
  case (rep,md) of (Dat rs, (_,ms)) ->
    map (get_site *** get_mod) (zip rs ms)
```

As this example shows, Forest blurs the distinction between data represented on disk and in memory. After writing a suitable Forest description, programmers can write programs that work on file system data as if it were in memory. Moreover, because Forest uses Haskell's lazy I/O operations, many simple programs do not require constructing an explicit representation of the entire directory being loaded in memory—a good thing as the directory of CoralCDN logs contains approximately 1GB of data! Instead, the load functions only read the portions of the file system that are needed to compute the result—in this case, only the site directories and not the gzipped log files contained within them.

As a final analysis example, consider a program that computes the top-*k* requested URLs from all CoralCDN nodes by size. The CoralCDN administrators compute this statistic periodically to help monitor and tune the performance of the system [11]. We define the analogous function in Haskell using helper functions such as `get_sites` to project out components of `rep`:

```
topk k =
  take k $ sortBy descBytes $ toList $
  fromListWith (+)
  [ (get_url e, get_total e)
  | (site,sdir) <- get_sites rep,
    (datetime,ldir) <- get_dates sdir,
    e <- get_entries ldir,
    is_in e ]
```

Reading this program inside-out, we see that it first uses a list comprehension to iterate through `rep`, collecting the individual log entries in the `coralwebsrv.log.gz` file for incoming requests and projecting out the URL requested and the total size of the request. It then sums the sizes of all requests for the same URL using the `fromListWith` function from the `Data.Map` module. Next, it sorts the entries in descending order. Finally, it returns the first *k* entries of the list as the final result.

Having implemented these analyses, a programmer may wish to store their results. She may do so via the following code, which uses `stats_manifest` to generate a manifest and `store` to copy it over to the `stats` directory. In addition, the code uses `stats_defaultMd`, a function that constructs default metadata for stats structures (a useful function in situations that require storing newly constructed data).

```
let result = Stats { last = sites_mod ()
                  , topk = topk 10 }
manifest <- stats_manifest
            ( result
            , stats_defaultMd result "/var/log/coral/stats" )
store manifest
```



Overall, the main take-away from this section is how Forest and its tight integration with Haskell facilitates exploratory data analysis, enabling remarkably terse queries over the combination of file contents, file attributes and directory structures.

## 5. Generic Tools

Third-party developers can use generic programming [18] to generate tools that will work for any filestore that has a Forest description. An advantage of these tools compared to tools that work directly on the untyped file system is that they are specific to the fragment of the file system relevant to the filestore. This fragment can be difficult to specify when using conventional tools since it can rely on the contents of configuration files, file naming conventions, file system attributes, *etc.* It is precisely these relationships that Forest descriptions capture concisely; tools written to use Forest specifications can leverage that information.

As a proof of concept, we have written a number of such tools, which we describe in this section.

### 5.1 Generic Querying

One simple application of generic programming is querying metadata to find files with a particular collection of attributes. The `findFiles` function

```
findFiles :: (ForestMD md) =>
  md -> (FileInfo -> Bool) -> [FilePath]
```

takes as input any Forest metadata value (*i.e.*, any value of type `md` where `md` belongs to the Forest metadata class `ForestMD`) and a predicate on `FileInfo` structures, and returns the list of all `FilePath`s anywhere in the input metadata whose associated `FileInfo` satisfies the predicate. For example, if `cs_md` is the metadata associated with the Princeton computer science department filestore, then the code

```
dirs = findFiles cs_md (\(r::FileInfo) ->
  (kind r) == DirectoryK)
other = findFiles cs_md (\(r::FileInfo) ->
  (owner r) /= "bwk")
```

binds `dirs` to the list of all directories in the data set and `other` to all the directories and files not owned by user "bwk".

To implement the `findFiles` function, we use the generic Haskell function `listify`:

```
findFiles md pred = map fullpath (listify pred md)
```

The return type of the polymorphic `listify` function is instantiated to match the argument type of its predicate argument. We map the `fullpath` function over the resulting list of `FileInfo` structures to return only the `FilePath`s.

### 5.2 File System Visualization

`ForestGraph` generates a graphical representation of any directory structure that matches a Forest specification. We generated the graphs in Figures 1 and 2 using this tool. In the default configuration, `ForestGraph` uses boxes to denote directories and ovals to denote files. Borders of varying thickness distinguish between ASCII and binary files. Dashed node boundaries indicate symbolic links and red nodes flag errors.

The core functionality of `ForestGraph` lies in the Haskell function `mdToPDF`:

```
mdToPDF :: ForestMD md =>
  md -> FilePath -> IO (Maybe String)
```

The function takes as input any metadata value and a filepath that specifies where to put the generated PDF file. It optionally returns a string (`Maybe String`); if the option is present, the string

contains an error message. The `IO` type constructor indicates that there can be side effects during the execution of the function. A use of this function to generate the graph for the Princeton computer science department filestore looks like:

```
do { (cs_rep,cs_md) <- CS_load "facadm"
  ; mdToPDF cs_md "Output/CS.pdf" }
```

Note that this code needs only the metadata to generate the graph; laziness means Forest will not load the representation in this case.

The related function `mdToPDFWithParams` takes an additional argument that allows the user to specify how to draw the nodes and edges in the output graph. Among other things, this parameter specifies how to map a value of type `Forest_md` into `GRAPHVIZ` [13, 14] attributes. By appropriately setting the parameter, a user can customize the formatting of each node according to its owner, group, or permissions, *etc.*, as well as specify global properties of the graph such as its orientation and size. `ForestGraph` uses the Haskell binding of the `GRAPHVIZ` library to lay out and render the graphs, so all customizations provided by `GRAPHVIZ` are available.

The `listify` function is at the heart of the implementation of this tool; we use it to convert the input metadata to the list of `FileInfos` in the metadata. We then convert this list into a graph data structure suitable for use with the `GRAPHVIZ` library.

### 5.3 Permission Checker

The permission tool is designed to check the permissions on the files and directories in a Forest description on a multi-user machine. In particular, it enables one user to determine which files a second user can read, write, or execute. If the second user cannot access a file in a particular way, the tool also reports the names of the files and directories whose permissions have to change to allow the access. The tool is useful when trying to share files with a colleague. It helps the first user ensure that all the necessary permissions have been set properly to allow the second user access. The key to the implementation of this tool is again applying the `listify` function to the metadata for the Forest description.

### 5.4 Shell Tools

We have implemented analogs of many shell tools that work over a file system fragment defined by a Forest description:

```
ls    :: (ForestMD md) => md -> String -> IO String
grep  :: (ForestMD md) => md -> String -> IO String
tar   :: (ForestMD md) => md -> FilePath -> IO ()
cp    :: (ForestMD md) => md -> FilePath -> IO ()
```

All of these functions work by extracting the relevant file names from the argument metadata structure using `listify` and then calling out to a shell tool to do the work. For `ls`, the second argument gives the command-line arguments to pass to the shell version of `ls`, and the result is the resulting output. The implementation uses `xarg` to lift the restriction on the number of files that can be passed to `ls`. For `grep`, the second argument is the search string and result is the output of the shell version of `grep`. For `tar`, the second argument specifies the location for the resulting tarball. The implementation uses a file manifest to allow `tar` to work regardless of the number of files involved. The `cp` tool uses the `tar` tool to move the files mentioned in the metadata to the location specified by the second argument *while retaining the same directory structure*. The module that implements these tools is 80 lines of Haskell code.

### 5.5 Description Inference Tool

This tool allows the user to generate a Forest description from the contents of the file system. The function

```
getDesc :: FilePath -> IO String
```

takes as an argument the path to the root of the directory structure to infer. It returns a string containing the generated representation. For example, below we show a fragment of the results when `getDesc` is invoked on the `classof11` directory:

```
data classof11 = Directory {
  aB11 is "AB11" :: aB11,
  bSE11 is "BSE11" :: bSE11,
  tTRANSFER is "TRANSFER" :: tTRANSFER,
  WITHDREW is "WITHDREW" :: WITHDREW }
data tTRANSFER = Directory {
  BEAUCHEMINtxt is "BEAUCHEMIN.txt" :: File Ptext,
  VERSTEEGtxt is "VERSTEEG.txt" :: File Ptext }
...
```

The description is not perfect: the label names are generated from the file name, for example. Nevertheless, the tool improves programmer productivity as it is easier for a programmer to edit a generated description than to start from scratch. Our first tool in this vein is simple; a more sophisticated variant would collapse records of files into lists when a width limit was exceeded or other criteria were met. Another variant might collapse deeply nested directories into a universal directory description when a depth limit was exceeded. The `getDesc` function works by using the universal description to load the contents of the file system starting from the supplied path. It then walks over the resulting metadata to generate a Forest parse tree, which it then pretty prints.

## 6. Implementation

The current implementation of Forest is available from the project web site: [forestproj.org](http://forestproj.org).

Haskell provides powerful language features and libraries that greatly facilitated implementation of Forest. The most obvious of these features is the quasiquotation mechanism [21] that we used to embed Forest into Haskell. This mechanism provided the benefits of being an embedded domain-specific language without having to sacrifice the flexibility of defining our own syntax. To use quasiquoting, we defined a Haskell value `forest` of type `QuasiQuoter` which specifies how to convert an input string representing a sequence of Forest declarations into the Template Haskell [27] data structures that represent the syntax of the corresponding collection of Haskell declarations. The Haskell compiler calls the `forest` “compilation” function during the compilation of any Haskell source file containing a Forest quasiquotation. The quasiquoted syntax `[forest| <input> |]` is legal anywhere the identifier `forest` is in scope. When the Haskell compiler processes this declaration, it first passes `<input>` as a string to the `forest` quasiquoter, and then it compiles the resulting Template Haskell data structures as if the corresponding Haskell code had appeared in the input at the location of the quasiquote. Early versions of quasiquoting supported quoting only expression and pattern forms. Simon Peyton Jones extended the mechanism to permit declaration and type quasi-quoting partly to enable the Forest implementation. We used this same approach to implement Pads/Haskell, which we built concomitantly.

Note that in implementing Forest, we had to use Template Haskell rather than any of the other libraries that support generic programming, both because that is what the quasiquote library expects and because we need to generate type and datatype declarations (and to do so at compile time). Other available generic libraries do not support the latter functionality.

**Parsing.** We used the parsec 3.1.0 parser combinator library [19] to implement the Forest parser. One key element of the Forest design is to allow arbitrary Haskell expressions in various places inside Forest descriptions. We did not want to reimplement the gram-

mar for Haskell expressions, which is quite complicated. Instead, we structured the Forest grammar so we could always determine the extent of any embedded Haskell code. We then used the Haskell Source Extension package [15] to parse these fragments. The data structure that this library returns is unfortunately not the data structure that Template Haskell requires, so we used yet another library, the Haskell Source Meta package [16], that provides this translation.

**Type checking.** We would like to give users high-quality error messages if there are type errors in their Forest declarations. At the moment, typechecking occurs, but only after the Forest declarations have been expanded to the corresponding Haskell code. Although these error messages can be quite informative, it is sub-optimal to report errors in terms of generated code. Type checking the Forest source is complicated by the embedded fragments of Haskell. As with the syntax, we do not want to reimplement the Haskell type-checker! There is an active proposal [29] to extend the Template Haskell infrastructure with functions that would enable us to ask the native Haskell typechecker for the types of embedded expressions and to extend the current type environment with type bindings for new identifiers. With this combination of features, we would be able to type check Forest sources directly.

**Runtime.** Although Forest facilitates treating the file system as a persistent store, it does not provide the ACID guarantees familiar from databases. None of the filestores we have encountered in practice are implemented in a system that provides such support; users instead have extra-linguistic mechanisms to make sure they do not corrupt their data with ill-timed concurrent reads and writes. That said, the Forest language does not preclude an implementation from providing such guarantees. We consider this issue very interesting future work.

Forest uses Haskell’s `unsafeInterleaveIO` to load each portion of a filestore only when needed by an application program. We have not measured the performance overhead of using Forest systematically. However, we have used our mostly-unoptimized implementation to manipulate filestores on the order of many gigabytes. Currently, the performance is acceptable for many applications.

The running time of storing operations is proportional to the “footprint” of the described filestore. However, the Forest compiler generates `load` and `manifest` functions for each named type in a description. Thus, updates may be made at any granularity for which there is a named type, which is typically at the level of individual files. We plan to investigate better support for incremental updates in future work.

## 7. A Core Calculus for Forest

This section describes a core calculus that formalizes the essential features of Forest precisely in a simple setting. It is inspired by classical (*i.e.*, not separating, substructural or ambient) unordered tree logics, customized for file systems. We used this calculus to investigate various features and prove theorems (such as the round-tripping properties presented at the end of this section) as we developed Forest.

### 7.1 The Basics: File Systems and Their Specifications

Figure 8 presents the formal file system model. A path  $r$  is a sequence of strings<sup>4</sup> and a file system  $F$  is finite map from paths to pairs of attributes  $a$  and file system contents  $T$ . We leave attributes abstract but expect that they include the usual fields: owner, group, date modified, *etc.* The attribute  $a_{\text{default}}$  contains default values for

<sup>4</sup>For simplicity, we ignore the special path elements “.” and “..”. It would be easy to add these features, at the cost of complicating the semantics.

### Basic definitions

<i>Integers</i>	$n \in \mathbb{Z}$
<i>Strings</i>	$u \in \Sigma^*$
<i>Booleans</i>	$b ::= \text{True} \mid \text{False}$
<i>Values</i>	$v ::= n \mid u \mid b \mid a \mid r \mid () \mid (v_1, v_2) \mid \text{Just } v \mid \text{Nothing} \mid [v_1, \dots, v_n]$
<i>Types</i>	$\tau ::= \text{Int} \mid \text{Bool} \mid () \mid (\tau_1, \tau_2) \mid \text{Maybe } \tau \mid [\tau]$
<i>Environments</i>	$\mathcal{E} ::= \emptyset \mid \mathcal{E}, x \mapsto v$
<i>Expressions</i>	$e ::= x \mid \lambda x. e \mid e_1 e_2 \mid \dots$

### Forest definitions

<i>Attributes</i>	$a ::= v$
<i>Paths</i>	$r ::= \cdot \mid r / u$
<i>Contents</i>	$T ::= \text{File } u \mid \text{Link } r \mid \text{Dir } \{u_1, \dots, u_n\}$
<i>File systems</i>	$F ::= \{r_1 \mapsto (a_1, T_1), \dots, r_n \mapsto (a_k, T_n)\}$
<i>Specifications</i>	$s ::= k_{\tau_1}^{\tau_2} \mid e :: s \mid \langle x : s_1, s_2 \rangle \mid [s \mid x : e] \mid \text{P}(e) \mid s ?$

**Figure 8.** Forest calculus syntax

all fields. The contents  $T$  of a node in the file system is either a file  $\text{File } n$  (where  $n$  is the string contents of the file), a symbolic link  $\text{Link } r$  (where  $r$  is the path pointed to by the link), or a directory  $\text{Dir } \{n_1, \dots, n_k\}$  (with paths  $n_1$  to  $n_k$ ). We write  $\text{dom}(F)$  for the set of paths  $F$  is defined on,  $F(r)$  for the contents at  $r$ , and  $F(r) = \perp$  when  $r$  is not in  $\text{dom}(F)$ .

A file system  $F$  is *well-formed* if it encodes a tree with directories at the internal nodes and files and symbolic links at the leaves. More formally,  $F$  is well-formed if the following conditions hold:

- $\text{dom}(F)$  is prefix-closed,
- $F(r) = (a, \text{Dir } \{n_1, \dots, n_k\}) \implies \forall i \in \{1, \dots, k\}. r / n_i \in \text{dom}(F)$ , and
- $F(r) = (a, \text{File } n_r) \vee F(r) = (a, \text{Link } r') \implies \forall n. r / n \notin \text{dom}(F)$ .

Note that although the structure of a well-formed file system is tree-shaped, cycles can be also expressed using symbolic links that point “upwards” in the file system.

Figure 8 also presents the syntax of file system specifications  $s$ . We leave the syntax of expression language abstract but assume that it contains values  $v$ , variables  $x$ , and other operators (of course, in the full Forest language, expressions can be arbitrary Haskell code). An environment  $\mathcal{E}$  maps variables to values. The semantic function  $\llbracket e \rrbracket_{\mathcal{E}, r}^{\tau}$  evaluates an expression  $e$  in the environment  $\mathcal{E}$  at path  $r$ , yielding a value  $v$  of type  $\tau$ .

These file system specifications are parameterized over a collection of constants  $k_{\tau_1}^{\tau_2}$ , which include specifications for files ( $\text{File}$ ), directories ( $\text{Dir}$ ), links ( $\text{Link}$ ), and Pads/Haskell-described files ( $\text{Adhoc}(b)$ ). The annotations  $\tau_1$  and  $\tau_2$  supply the internal types of the the representation and the constant-specific portion of the meta-data. The meta-variable  $b$  ranges over bidirectional functions: in the forward direction, such functions load (parse) data; in the reverse direction they store (print) it.

To define the semantics of the overall language precisely, we assume that each constant is associated with functions  $\text{load}_k$  and  $\text{store}_k$ . For example, the  $\text{load}$  function for the  $\text{File}$  construct, which describes any file (but not symbolic links or directories), is defined

as follows:

$$\text{load}_{\text{File}}(\mathcal{E}, F, r) = \begin{cases} (n, (\text{True}, a)), & \text{if } F(r) = (a, \text{File } n) \\ ("", (\text{False}, a_{\text{default}})), & \text{otherwise} \end{cases}$$

The arguments to the function include an environment  $\mathcal{E}$ , a file system  $F$  to load from, and a path  $r$  within that file system. This function either returns the contents and attributes of the file at path  $r$ , if it exists, or “” and default attributes if  $F$  does not contain a file at  $r$ . The  $\text{store}$  function for  $\text{File}$  is defined as follows:

$$\text{store}_{\text{File}}(\mathcal{E}, F, r, v, d) = \begin{cases} (F[r := (a, \text{File } v)], & \text{if } d = (\text{True}, a) \\ \phi' = \lambda F'. (F'(r) = (a, \text{File } v))) & \\ F[r := \perp], & \text{if } d = (\text{False}, a) \wedge \\ \phi' = \lambda F'. F'(r) \neq (\_, \text{File } \_) & F(r) = (\_, \text{File } \_) \\ F, & \text{otherwise} \\ \phi' = \lambda F'. F'(r) \neq (\_, \text{File } \_) & \end{cases}$$

The arguments to the store function include an environment  $\mathcal{E}$ , an existing file system  $F$  to store into, a path  $r$  to store at, a representation  $v$  for the data to store, and the metadata  $d$  associated with that representation. The store function produces two results: an updated file system  $F'$  and a predicate  $\phi'$  that records the constraints needed to ensure consistency. In this case, the  $\text{store}$  function for  $\text{File}$  overwrites the contents of the file system  $F$  at path  $r$  with  $(a, \text{File } v)$  if  $d$  is valid (and contains  $a$ ), deletes the contents of  $F$  at  $r$  if  $d$  is not valid but  $F(r)$  contains a file, and otherwise returns  $F$  unchanged. The predicate  $\phi'$  requires that  $F'(r)$ , the contents of the new file system  $F'$  at  $r$ , be  $\text{File } v$  in the first case and that  $F'(r)$  not be a file in the other two cases. These constraints must be satisfied in order to guarantee the round-tripping properties presented at the end of this section.

In the Forest surface syntax, records and paths are specified using a single construct (and similarly for comprehensions) while the core calculus models (dependent) records, paths, and comprehensions as independent, orthogonal constructs. Path specifications are written  $e :: s$ , where  $e$  is a path name (to be appended to the current path) and  $s$  specifies a fragment of the file system at that path. Record specifications are written  $\langle x : s_1, s_2 \rangle$ , where  $x$  may appear in  $s_2$ . Comprehension specifications are written  $[s \mid x : e]$ , where  $e$  is an expression that describes a set of values,  $x$  is a variable, and  $s$ , which may depend on  $x$ , specifies a fragment of the file system for each value of  $x$ . For example, the specification

```
{c is "c.txt" :: C, d is "d.txt" :: D c}
```

is encoded in the calculus as  $\langle x : ("c.txt" :: C), ("d.txt" :: D x) \rangle$ . Similarly, the comprehension

```
[c :: C | c <- matches (GL "*" )]
```

is encoded as  $\langle x : \text{Dir}, [y :: C \mid y : x] \rangle$ . Predicate specifications  $\text{P}(e)$  succeed when  $e$  evaluates to  $\text{True}$  and fail when  $e$  evaluates to  $\text{False}$  under the current environment. A Forest constraint of the form  $s$  **where**  $e$  is encoded in the calculus using a dependent pair and a predicate:  $\langle x : s, \text{P}(e[x/\text{this}]) \rangle$ . Finally, maybe specifications are written as  $s ?$  in the calculus.

## 7.2 Calculus Semantics

The semantics of the calculus is organized into four separate definitions, one for each of the four major artifacts generated by the Forest compiler.

**Type Definitions.** Figure 10 defines types for the representations  $\mathcal{R}[\![s]\!]$  and metadata  $\mathcal{M}[\![s]\!]$  generated by specifications  $s$ . The types for constants  $k_{\tau_1}^{\tau_2}$  are read off from their annotations while the

$$\boxed{\mathcal{E}; r; s \vdash \text{load } F \triangleright (v, d)}$$

$$\frac{}{\mathcal{E}; r; k_{\tau_1}^{T_2} \vdash \text{load } F \triangleright (\text{load}_k(\mathcal{E}, F, r))}$$

$$\frac{\mathcal{E}; \llbracket r/e \rrbracket_{\text{path}}^{\mathcal{E}, r}; s \vdash \text{load } F \triangleright (v, d)}{\mathcal{E}; r; e::s \vdash \text{load } F \triangleright (v, d)}$$

$$\frac{\mathcal{E}; r; s_1 \vdash \text{load } F \triangleright (v_1, d_1) \quad (\mathcal{E}, x_{\text{rep}} \mapsto v_1, x_{\text{md}} \mapsto d_1); r; s_2 \vdash \text{load } F \triangleright (v_2, d_2) \quad b = \text{valid}(d_1) \wedge \text{valid}(d_2)}{\mathcal{E}; r; \langle x:s_1, s_2 \rangle \vdash \text{load } F \triangleright ((v_1, v_2), (b, (d_1, d_2)))}$$

$$\frac{\begin{array}{l} \llbracket e \rrbracket_{\text{tau}}^{\mathcal{E}, r} = [w_1, \dots, w_k] \\ \forall i \in \{1, \dots, k\}. (\mathcal{E}, x \mapsto w_i); r; s \vdash \text{load } F \triangleright (v_i, d_i) \\ b = \bigwedge_i^k \text{valid}(d_i) \quad vs = [v_1, \dots, v_k] \quad ds = [d_1, \dots, d_k] \end{array}}{\mathcal{E}; r; [s \mid x : e] \vdash \text{load } F \triangleright (vs, (b, ds))}$$

$$\frac{b = \llbracket e \rrbracket_{\text{bool}}^{\mathcal{E}, r}}{\mathcal{E}; r; \text{P}(e) \vdash \text{load } F \triangleright ((\text{), (b, (\text{))}))}$$

$$\frac{r \notin \text{dom}(F)}{\mathcal{E}; r; s_1? \vdash \text{load } F \triangleright (\text{Nothing}, (\text{True}, \text{Nothing}))}$$

$$\frac{r \in \text{dom}(F) \quad \mathcal{E}; r; s_1 \vdash \text{load } F \triangleright (v_1, d_1)}{\mathcal{E}; r; s_1? \vdash \text{load } F \triangleright (\text{Just } v_1, (\text{valid}(d_1), \text{Just } d_1))}$$

(a)

$$\boxed{\mathcal{E}; r; s \vdash \text{store } (F, v, d) \triangleright (F', \phi')}$$

$$\frac{}{\mathcal{E}; r; k_{\tau_1}^{T_2} \vdash \text{store } (F, v, d) \triangleright (\text{store}_k(\mathcal{E}, F, r, v, d))}$$

$$\frac{\mathcal{E}; \llbracket r/e \rrbracket_{\text{path}}^{\mathcal{E}, r}; s \vdash \text{store } (F, v, d) \triangleright (F', \phi')}{\mathcal{E}; r; e::s \vdash \text{store } (F, v, d) \triangleright (F', \phi')}$$

$$\frac{\mathcal{E}; r; s_1 \vdash \text{store } (F, v_1, d_1) \triangleright (F'_1, \phi'_1) \quad (\mathcal{E}, x_{\text{rep}} \mapsto v_1, x_{\text{md}} \mapsto d_1); r; s_2 \vdash \text{store } (F, v_2, d_2) \triangleright (F'_2, \phi'_2) \quad \phi' = \lambda F'. (b = \text{valid}(d_1) \wedge \text{valid}(d_2)) \wedge \phi'_1(F') \wedge \phi'_2(F')}{\mathcal{E}; r; \langle x:s_1, s_2 \rangle \vdash \text{store } (F, (v_1, v_2), (b, (d_1, d_2))) \triangleright (F'_1++F'_2, \phi')}$$

$$\frac{\begin{array}{l} vs = [v_1, \dots, v_j] \quad ds = [d_1, \dots, d_l] \\ \llbracket e \rrbracket_{\text{tau}}^{\mathcal{E}, r} = [w_1, \dots, w_m] \quad k = \min(j, l, m) \\ \forall i \in \{1, \dots, k\}. (\mathcal{E}, x \mapsto w_i); r; s \vdash \text{store } (F, v_i, d_i) \triangleright (F'_i, \phi'_i) \\ \phi' = \lambda F'. (j = l = m) \wedge (b = \bigwedge_i^k \text{valid}(d_i)) \wedge (\bigwedge_i^k \phi'_i(F')) \end{array}}{\mathcal{E}; r; [s \mid x : e] \vdash \text{store } (F, vs, (b, ds)) \triangleright (F'_1++\dots++F'_k, \phi')}$$

$$\frac{\phi' = \lambda F'. (b = \llbracket e \rrbracket_{\text{bool}}^{\mathcal{E}, r})}{\mathcal{E}; r; \text{P}(e) \vdash \text{store } (F, (\text{), (b, (\text{)))} \triangleright (F', \phi')}$$

$$\frac{\mathcal{E}; r; s_1 \vdash \text{store } (F, v_1, d_1) \triangleright (F', \phi'_1) \quad \phi' = \lambda F'. (b = \text{valid}(d_1)) \wedge (r \in \text{dom}(F')) \wedge \phi'_1(F')}{\mathcal{E}; r; s_1? \vdash \text{store } (F, \text{Just } v_1, (b, \text{Just } d_1)) \triangleright (F', \phi')}$$

$$\frac{\phi' = \lambda F'. (d = \text{Nothing}) \wedge b \wedge r \notin \text{dom}(F')}{\mathcal{E}; r; s_1? \vdash \text{store } (F, \text{Nothing}, (b, d)) \triangleright (F[r := \perp], \phi')}$$

$$\frac{\mathcal{E}; r; s_1 \vdash \text{store } (F, v_1, d_{\text{default}}^{s_1}) \triangleright (F', \phi'_1) \quad \phi' = \lambda F'. \text{False}}{\mathcal{E}; r; s_1? \vdash \text{store } (F, \text{Just } v_1, (b, \text{Nothing})) \triangleright (F', \phi')}$$

(b)

Figure 10. Forest calculus semantics for (a) loading and (b) storing

$s$	$\mathcal{R}[s]$	$\mathcal{M}[s]$
$k_{\tau_1}^{T_2}$	$\tau_1$	$Md \tau_2$
$e::s$	$\mathcal{R}[s]$	$\mathcal{M}[s]$
$\langle x:s_1, s_2 \rangle$	$(\mathcal{R}[s_1], \mathcal{R}[s_2])$	$Md (\mathcal{M}[s_1], \mathcal{M}[s_2])$
$[s \mid x : e]$	$\mathcal{R}[s]$	$Md [\mathcal{M}[s]]$
$\text{P}(e)$	$()$	$Md ()$
$s?$	$\text{Maybe } \mathcal{R}[s]$	$Md (\text{Maybe } \mathcal{M}[s])$

Figure 9. Forest calculus representation and metadata types

types for other specifications are constructed from their structure in the obvious way—*e.g.*, the type of representations for  $\langle x:s_1, s_2 \rangle$  is a product  $(\mathcal{R}[s_1], \mathcal{R}[s_2])$ . The type constructor  $Md$  provides a uniform representation for metadata and is defined as follows:

$$\begin{aligned} Md \tau &= (\text{Header}, \tau) \\ \text{Header} &= \text{Bool} \end{aligned}$$

The function  $\text{valid}(d)$  extracts the boolean from the metadata structure  $d$ , returning *True* if there are no errors in the structure and *False* otherwise.

**Semantics of Loading and Storing.** The inference rules on the left side of Figure 10 define the semantics of the load function. Reading from right to left, the judgment  $\mathcal{E}; r; s \vdash \text{load } F \triangleright (v, d)$  states one can obtain the pair  $(v, d)$  of representation and metadata, by materializing components of the filesystem  $F$  in memory using the specification  $s$  at path  $r$  in environment  $\mathcal{E}$ . Reading from left to right, this judgment may also be viewed as a total function from  $\mathcal{E}, r, s$  and  $F$  to  $(v, d)$ . The judgment is total because when  $F$  fails to match  $s$ , the load function generates defaults in the representation  $v$  and records errors in the metadata  $d$ . This design allows a programmer to explore a file system fragment even when it does not match the given specification exactly.

Let us examine a few of the inference rules that define the store function in detail. The rule for constants  $k_{\tau_1}^{T_2}$  just invokes the associated  $\text{store}_k$  function. The rule for  $[s \mid x : e]$  comprehensions is more interesting: it first evaluates  $e$  to a list  $[w_1, \dots, w_k]$  and then invokes the store function for  $s$   $k$  times in environments where  $x$  is bound to each  $w_i$ . It then collects up the results into lists

of representations  $[v_1, \dots, v_k]$  and metadata  $[d_1, \dots, d_k]$ , which it uses as the final result. The predicate  $P(e)$  construct tests whether an expression  $e$  is satisfied. It returns  $()$  as the representation and  $([e]_{bool}^{\mathcal{E}, r}, ())$  as the metadata. Finally,  $s?$  invokes  $s$ 's load function if the current path  $r$  exists in the file system, injecting the result into the maybe type using `Just`, and otherwise returns `Nothing`.

The inference rules on the right side of Figure 10 define the store function. The judgment  $\mathcal{E}; r; s \vdash \text{store } (F, v, d) \triangleright (F', \phi')$  states that in environment  $\mathcal{E}$  storing  $(v, d)$  into file system  $F$  using specification  $s$  yields the file system  $F'$  and predicate  $\phi'$ . The predicate  $\phi'$  tracks the conditions on the file system needed to ensure that it accurately reflects the information in the representation and metadata.

As a simple example to illustrate why predicates are needed, consider the specification  $s = \langle x:\text{File}, \text{File} \rangle$  and suppose that the load function is called in an environment  $\mathcal{E}$  with a file system  $F$  and path  $r$  where  $F(r) = (a, \text{File } n)$ . The representation returned by `load` will be a pair  $(n, n)$  containing two copies of the file contents at  $r$  and the metadata will also contain a pair  $(\text{True}, (\text{True}, a), (\text{True}, a))$  with two copies of the metadata associated with that file. Now suppose that we change the representation to  $(n, n')$ , with  $n \neq n'$ , and we store the result back to the file system. Unfortunately, because the representation is inconsistent—it does not satisfy the dependency between the two components of the pair implied by  $s$ —the store function cannot produce a new file system containing the information in both  $n$  and  $n'$ . Thus, it must store one and discard the other. The predicate  $\phi'$  generated by the store function provides a way to track and report inconsistencies. In this case, the predicate will be equivalent to the following:

$$\phi' = \lambda F'. (F'(r) = (a, \text{File } n)) \wedge (F'(r) = (a, \text{File } n'))$$

which is obviously not satisfiable when  $n \neq n'$ .

Now that the overall structure of the store judgement has been explained, let us examine a few of the inference rules in detail. The rule for constants  $k_{\tau_1}^{\tau_2}$  simply invokes the  $\text{store}_k$  function. The rule for path specifications  $e::s$  passes off control to the store function for  $s$  after replacing the current path  $r$  with  $\llbracket r/e \rrbracket_{path}^{\mathcal{E}, r}$ . The rule for dependent pairs  $\langle x:s_1, s_2 \rangle$  is more interesting. Given a pair  $(v_1, v_2)$  as the representation, it first invokes the store function for  $s_1$  with  $v_1$ , producing an updated file system  $F'_1$  and predicate  $\phi'_1$ . Next, it invokes the store function for  $s_2$  with  $v_2$  in an extended environment where  $x$  is bound to  $v_1$ , yielding another updated file system  $F'_2$  and  $\phi'_2$ . It combines the updated file systems using the following right-biased concatenation operator,

$$(F_1 ++ F_2)(r) = \begin{cases} (a_2, \text{Dir } N_1 \cup N_2) & \text{if } F_1(r) = (a_1, \text{Dir } N_1) \wedge \\ & F_2(r) = (a_2, \text{Dir } N_2) \\ F_1(r) & \text{if } F_2(r) = \perp \\ F_2(r) & \text{otherwise} \end{cases}$$

Finally, it combines the predicates using conjunction. The result is a file system that contains the consistent changes made to the file system by the store functions for  $s_1$  and  $s_2$  as well as a predicate that checks for the consistency of all of their changes.

### 7.3 Formal Properties

The first property of the Forest calculus is a basic type safety property, which states that the load function for specifications  $s$  generates representations and metadata belonging to  $\mathcal{R}[\![s]\!]$  and  $\mathcal{M}[\![s]\!]$  respectively.

#### Proposition 1 (Type Safety)

If  $\mathcal{E}; r; s \vdash \text{load } F \triangleright (v, d)$  and  $\mathcal{R}[\![s]\!] = \tau_{\mathcal{R}}$  and  $\mathcal{M}[\![s]\!] = \tau_{\mathcal{M}}$  then  $\vdash v : \tau_{\mathcal{R}}$  and  $\vdash d : \tau_{\mathcal{M}}$ .

The above property demonstrates that our type definitions are properly aligned with the semantics of loading. To ensure that the semantics of loading is, in turn, aligned with the semantics of storing, we also prove the following two round-tripping properties.

#### Theorem 2 (LoadStore)

Let  $\mathcal{E}$  be an environment,  $F$  a file system,  $r$  a path,  $s$  a specification,  $v$  a representation, and  $d$  metadata. If

$$\begin{aligned} \mathcal{E}; r; s \vdash \text{load } F \triangleright (v, d) \\ \mathcal{E}; r; s \vdash \text{store } (F, v, d) \triangleright (F', \phi) \end{aligned}$$

then  $F = F'$  and  $\phi(F')$ .

#### Theorem 3 (StoreLoad)

Let  $\mathcal{E}$  be an environment,  $F$  and  $F'$  file systems,  $r$  a path,  $s$  a specification,  $v$  a representation,  $d$  and  $d'$  metadata, and  $\phi'$  a predicate. If

$$\begin{aligned} \mathcal{E}; r; s \vdash \text{store } (F, v, d) \triangleright (F', \phi') \\ \phi'(F') \\ \mathcal{E}; r; s \vdash \text{load } F' \triangleright (v', d') \end{aligned}$$

then  $v' = v$  and  $\text{valid}(d) = \text{valid}(d')$ .

The first theorem states that loading from a file system  $F$  and immediately storing the resulting representation and metadata yields the original file system and, moreover, it satisfies the predicate produced by the store function. The second theorem states that storing an arbitrary representation and metadata and then loading the resulting file system yields the same representation and contains errors only if the original metadata also contained errors. These properties are based on the general correctness conditions that have been proposed for bidirectional transformations in the context of lenses [10], but are generalized here to accommodate the inconsistencies that can arise when working with imperfect, ad hoc data. The proofs of these theorems can be found in Appendix A.

## 8. Related Work

The work in this paper builds upon ideas developed in the Pads project [5, 7]. Pads uses extended type declarations to describe the grammar of a document and simultaneously to generate types for parsed data and a suite of data-processing tools. The obvious difference between Pads (and other parser generators) and Forest is that Pads generates infrastructure for processing strings (the insides of a single file) whereas Forest generates infrastructure for processing entire file systems. Forest (and Pads/Haskell) is architecturally superior to previous versions of Pads in the tight integration with its host language and in its support for third-party generic programming and tool construction.

More generally, Forest shares high-level goals with other systems that seek to make data-oriented programming simpler and more productive. For example, Microsoft's LINQ [20] extends the .NET languages to enable querying any data source that supports the `IEnumerable` interface using a simple, convenient syntax. LINQ differs in that it does not provide support for declaratively specifying the structure of, and then ingesting, filestores. *Type Providers* [28], an experimental feature of F#, help programmers materialize standard data sources equipped with predefined schemas (such as XML documents or databases) in memory in an F# program. Type Providers do not themselves provide a new means for describing data sources (as Forest does).

Several XML-based languages for specifying file formats, file organization and file locations have been proposed. One example of such a language is XFiles [1]. XFiles uses RDF specifications to describe the location, permissions, ownership, and other attributes of

files, as well as the name of an application capable of parsing specific files. The key difference between XFiles and Forest is that Forest is tightly integrated into a general-purpose, conventional programming language. Forest declarations generate types, functions and data structures that materialize the data within a surrounding Haskell program while XFiles does not interoperate directly with a conventional programming language.

A recent MSc thesis by Ntzik proposes using an extension of context logic [2] to reason about the effects of updates made to file systems using standard POSIX commands [24]. The core goal of Ntzik's work is to create a new kind of Hoare Logic, and consequently, it is quite different from Forest. In addition, technically, Forest is more closely related to classical tree logics than to substructural logics such as context logic.

The round-tripping properties that core Forest programs obey are based on laws that have been proposed in the context of well-behaved bidirectional transformations, often called lenses [10]. As far as we are aware, lenses for file systems have not been developed but some of the same fundamental issues that arise in core Forest have been studied by Hu and his colleagues, including handling data with internal dependencies [23] as well as graph structures [17].

## 9. Conclusions

In this paper, we present the design of Forest, an embedded domain-specific language for describing filestores. A Forest description concisely specifies a collection of files, directories, and symbolic links as well as expected file system attributes such as owners and permissions. From a description, the Forest compiler generates code to lazily load the on-disk data into an isomorphic in-memory representation, lowering the divide between on-disk and in-memory data. Forest also generates type class instances that make it easy for third-party tool developers to use Haskell's generic programming infrastructure. We have used this infrastructure ourselves to define a number of useful tools. In addition, the language has a formal semantics based on classical tree logics and is fully implemented. On the latter point, our work serves as an extensive case study in domain-specific language design, and, as such, has inspired changes in the design of Template Haskell. Source code for Forest is available from the Forest web site [9].

## Acknowledgments

We would like to thank Simon Peyton Jones for extending Haskell's Quasi-Quoting and Template Haskell mechanisms to support the Forest design and for help in using the mechanisms in general.

This material is based upon work supported under NSF grant CCF-1016937 and ONR grant N00014-09-1-0652. Any opinions, findings, and recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the NSF or the ONR. Kenny Q. Zhu was partially supported by NSFC Grant No. 61033002.

## References

[1] S.-C. Buraga. An XML-based semantic description of distributed file systems. In *RoEduNet International Conference on Networking in Education and Research*, pages 41–48, 2003.

[2] C. Calcagno, P. Gardner, and U. Zarfaty. Context logic and tree update. In *POPL*, pages 271–282, 2005.

[3] Filesystem Hierarchy Standard Group. Filesystem hierarchy standard. <http://www.pathname.com/fhs/>, 2004.

[4] K. Fisher, N. Foster, D. Walker, and K. Q. Zhu. Forest: A Language and Toolkit for Programming with Filestores. Technical Report., June 2011.

[5] K. Fisher and R. Gruber. PADS: A domain specific language for processing ad hoc data. In *PLDI*, pages 295–304, June 2005.

[6] K. Fisher, Y. Mandelbaum, and D. Walker. The next 700 data description languages. In *POPL*, Jan. 2006.

[7] K. Fisher, Y. Mandelbaum, and D. Walker. The next 700 data description languages. *JACM*, 57:10:1–10:51, February 2010.

[8] K. Fisher and D. Walker. The PADS project: An overview. In *Proceedings of the 14th International Conference on Database Theory, ICDT '11*, pages 11–17, New York, NY, USA, 2011. ACM.

[9] Forest: A language and toolkit for programming with file system fragments. <http://forestproj.org>, 2010.

[10] J. N. Foster, M. B. Greenwald, J. T. Moore, B. C. Pierce, and A. Schmitt. Combinators for bidirectional tree transformations: A linguistic approach to the view update problem. *TOPLAS*, 29(3), May 2007.

[11] M. J. Freedman. Experiences with CoralCDN: A five-year operational view. In *NSDI*, pages 7–7, 2010.

[12] M. J. Freedman, E. Freudenthal, and D. Mazieres. Democratizing content publication with Coral. In *NSDI*, pages 18–18, 2004. See also <http://www.coralcdn.org/>.

[13] E. R. Gansner and S. C. North. An open graph visualization system and its applications to software engineering. *Softw. Pract. Exper.*, 30:1203–1233, September 2000.

[14] Haskell Graphviz Package. <http://hackage.haskell.org/package/graphviz>.

[15] Haskell Source Extensions Package. <http://hackage.haskell.org/package/haskell-src-exts>.

[16] Haskell Source Meta Package. <http://hackage.haskell.org/package/haskell-src-meta>.

[17] S. Hidaka, Z. Hu, K. Inaba, H. Kato, K. Matsuda, and K. Nakano. Bidirectionalizing graph transformations. In *ICFP*, pages 205–216, 2010.

[18] R. Lämmel and S. P. Jones. Scrap your boilerplate: A practical design pattern for generic programming. In *TLDI*, pages 26–37, 2003.

[19] D. Leijen and E. Meijer. Parsec: Direct style monadic parser combinators for the real world. Technical Report UU-CS-2001-27, Department of Computer Science, Universiteit Utrecht, 2001.

[20] LINQ: .NET language-integrated query. <http://msdn.microsoft.com/library/bb308959.aspx>, Feb. 2007.

[21] G. Mainland. Why it's nice to be quoted: Quasiquoting for Haskell. In *Haskell Workshop*, pages 73–82, 2007.

[22] Y. Mandelbaum, K. Fisher, D. Walker, M. Fernández, and A. Gleyzer. PADS/ML: A functional data description language. In *POPL*, Jan. 2007.

[23] S.-C. Mu, Z. Hu, and M. Takeichi. An algebraic approach to bidirectional updating. In *ASIAN Symposium on Programming Languages and Systems (APLAS)*, pages 2–20, Nov. 2004.

[24] G. Ntzik. Local reasoning for filesystems. Master's thesis, Imperial College, Sept. 2010.

[25] PADS project. <http://www.padsproj.org/>, 2007.

[26] T. J. Parr and R. W. Quong. ANTLR: A predicated- $LL(k)$  parser generator. *Software—Practice and Experience*, 25(7):789–810, July 1995.

[27] T. Sheard and S. P. Jones. Template meta-programming for Haskell. In *Haskell Workshop*, pages 1–16, 2002.

[28] D. Syme. Looking Ahead with F#: Taming the Data Deluge. Presentation at the Workshop on F# in Education, Nov. 2010.

[29] Template Haskell Extension Proposal. [hackage.haskell.org/trac/ghc/blog/Template%20Haskell%20Proposal](http://hackage.haskell.org/trac/ghc/blog/Template%20Haskell%20Proposal).

## A. Proofs

This appendix contains the proofs of the theorems stated in Section 7.

### Theorem 2 (LoadStore)

Let  $\mathcal{E}$  be an environment,  $F$  a file system,  $r$  a path,  $s$  a specification,  $v$  a representation, and  $d$  metadata. If

$$\begin{aligned} \mathcal{E}; r; s \vdash \text{load } F \triangleright (v, d) \\ \mathcal{E}; r; s \vdash \text{store } (F, v, d) \triangleright (F', \phi) \end{aligned}$$

then  $F = F'$  and  $\phi'(F')$ .

**Proof:** The proof is by induction on  $s$ . For every constant  $k_{\tau_1}^2$ , we assume that  $\text{load}_k$  and  $\text{store}_k$  satisfy the theorem.

**Case:**  $s = k_{\tau_1}^2$

By the definitions of the load and store functions we have

$$\begin{aligned} v, d &= \text{load}_k(\mathcal{E}, F, r) \\ F', \phi' &= \text{store}_k(\mathcal{E}, F, r, v, d) \end{aligned}$$

By the assumptions about the behavior of  $\text{load}_k$  and  $\text{store}_k$ , we have  $F' = F$  and  $\phi'(F')$ , which finishes the case.

**Case:**  $s = e::s_1$

By the definitions of the load and store functions we have

$$\begin{aligned} \mathcal{E}; r'; s_1 \vdash \text{load } F \triangleright (v, d) & & r' &= \llbracket r / e \rrbracket_{path}^{\mathcal{E}, r} \\ \mathcal{E}; r'; s_1 \vdash \text{store } (F, v, d) \triangleright (F', \phi') & & & \end{aligned}$$

By the induction hypothesis applied to  $s_1$  we have  $F' = F$  and  $\phi'(F')$ , which finishes the case.

**Case:**  $s = \langle x:s_1, s_2 \rangle$

By the definitions of the load and store functions we have

$$\begin{aligned} \mathcal{E}; r; s_1 \vdash \text{load } F \triangleright (v_1, d_1) & & \mathcal{E}' &= \mathcal{E}, x_{rep} \mapsto v_1, x_{md} \mapsto d_1 \\ \mathcal{E}'; r; s_2 \vdash \text{load } F \triangleright (v_2, d_2) & & b &= \text{valid}(d_1) \wedge \text{valid}(d_2) \\ \mathcal{E}; r; s_1 \vdash \text{store } (F, v_1, d_1) \triangleright (F'_1, \phi'_1) & & v &= (v_1, v_2) \\ \mathcal{E}'; r; s_2 \vdash \text{store } (F, v_2, d_2) \triangleright (F'_2, \phi'_2) & & d &= (b, (d_1, d_2)) \\ & & F' &= F'_1 ++ F'_2 \\ & & \phi' &= \lambda F'. (b = \text{valid}(d_1) \wedge \text{valid}(d_2)) \wedge \phi'_1(F') \wedge \phi'_2(F') \end{aligned}$$

By the induction hypothesis applied to  $s_1$ , we have  $F'_1 = F$  and  $\phi'_1(F'_1)$ . Likewise, by the induction hypothesis applied to  $s_2$ , we have  $F'_2 = F$  and  $\phi'_2(F'_2)$ . We immediately have  $\phi'(F')$  and  $F' = F$  as  $(++)$  is idempotent, which finishes the case.

**Case:**  $s = [s_1 \mid x : e]$

By the definitions of the load and store functions we have

$$\begin{aligned} \llbracket e \rrbracket_{[tau]}^{\mathcal{E}, r} &= [w_1, \dots, w_k] & v &= [v_1, \dots, v_k] \\ \forall i \in \{1, \dots, k\}. (\mathcal{E}, x \mapsto w_i); r; s_1 \vdash \text{load } F \triangleright (v_i, d_i) & & b &= \bigwedge_i^k \text{valid}(d_i) \\ \forall i \in \{1, \dots, k\}. (\mathcal{E}, x \mapsto w_i); r; s \vdash \text{store } (F, v_i, d_i) \triangleright (F'_i, \phi'_i) & & d &= (b, [d_1, \dots, d_k]) \\ & & F' &= F'_1 ++ \dots ++ F'_k \\ & & \phi' &= \lambda F'. (b = \bigwedge_i^k \text{valid}(d_i)) \wedge (\bigwedge_i^k \phi'_i(F')) \end{aligned}$$

By the induction hypothesis applied to  $s$  ( $k$  times), we have  $F'_i = F$  and  $\phi'_i(F'_i)$  for  $i$  from 1 to  $k$ . We immediately have  $\phi'(F')$  and  $F' = F$  as  $(++)$  is idempotent, which finishes the case.

**Case:**  $s = P(e)$

By the definition of the load and store functions we have

$$\begin{aligned} v &= () & F' &= F \\ d &= (b, ()) & \phi' &= \lambda F'. (b = \llbracket e \rrbracket_{bool}^{\mathcal{E}, r}) \\ b &= \llbracket e \rrbracket_{bool}^{\mathcal{E}, r} & & \end{aligned}$$

Thus, we immediately have  $F' = F$  and  $\phi'(F')$ .

**Case:**  $s = s_1?$

We analyze two subcases:

**Subcase:**  $r \in \text{dom}(F)$

By the definition of the load and store functions we have

$$\begin{aligned} \mathcal{E}; r; s \vdash \text{load } F \triangleright (v_1, d_1) & & v &= \text{Just } v_1 \\ \mathcal{E}; r; s \vdash \text{store } (F, v_1, d_1) \triangleright (F', \phi'_1) & & d &= (b, \text{Just } d_1) \\ & & b &= \text{valid}(d_1) \\ & & \phi' &= \lambda F'. (b = \text{valid}(d_1)) \wedge (r \in \text{dom}(F)) \wedge \phi'_1(F') \end{aligned}$$

By the induction hypothesis applied to  $s_1$  we have  $F' = F$  and  $\phi'_1(F')$ . We immediately have  $\phi'(F')$ .

**Subcase:**  $r \notin \text{dom}(F)$

By the definition of the load and store functions we have

$$\begin{aligned} v &= \text{Nothing} \\ d &= (b, d_1) & F' &= F[r := \perp] \\ d_1 &= \text{Nothing} & \phi' &= \lambda F'. (d_1 = \text{Nothing}) \wedge b \wedge (r \notin \text{dom}(F')) \\ b &= \text{True} \end{aligned}$$

We immediately have  $F' = F$  and  $\phi'(F')$ , which finishes the case and the inductive proof. ■

### Theorem 3 (StoreLoad)

Let  $\mathcal{E}$  be an environment,  $F$  and  $F'$  file systems,  $r$  a path,  $s$  a specification,  $v$  a representation,  $d$  and  $d'$  metadata, and  $\phi'$  a predicate. If

$$\begin{aligned} \mathcal{E}; r; s \vdash \text{store } (F, v, d) \triangleright (F', \phi') \\ \mathcal{E}; r; s \vdash \text{load } F' \triangleright (v', d') \end{aligned}$$

then  $v' = v$  and  $\text{valid}(d) = \text{valid}(d')$ .

**Proof:** We will prove a slightly stronger result that implies the theorem: for all environments  $\mathcal{E}$ , file systems  $F$ ,  $G_1$ ,  $G_2$ , and  $F'$ , paths  $r$ , specifications  $s$ , representations  $v$  and  $v'$ , metadata  $d$  and  $d'$ , and constraints  $\phi'$ , if

$$\begin{aligned} \mathcal{E}; r; s \vdash \text{store } (F, v, d) \triangleright (F', \phi') \\ \mathcal{E}; r; s \vdash \text{load } (G_1 ++ F' ++ G_2) \triangleright (v', d') \end{aligned}$$

then  $v' = v$  and  $\text{valid}(d) = \text{valid}(d')$ .

The proof is by induction on  $s$ . For every constant  $k_{\tau_1}^{\tau_2}$ , we assume that  $\text{load}_k$  and  $\text{store}_k$  satisfy the strengthened property.

**Case :**  $s = k_{\tau_1}^{\tau_2}$

By the definitions of the load and store functions we have

$$\begin{aligned} F', \phi &= \text{store}_k(\mathcal{E}, F, r, v, d) \\ v', d' &= \text{load}_k(\mathcal{E}, (G_1 ++ F' ++ G_2), r s) \end{aligned}$$

By assumptions about the behavior of  $\text{load}_k$  and  $\text{store}_k$ , we have  $v' = v$  and  $\text{valid}(d) = \text{valid}(d')$ , which finishes the case.

**Case:**  $s = e :: s_1$

By the definitions of the load and store functions we have

$$\begin{aligned} \mathcal{E}; r'; s_1 \vdash \text{store } (F, v, d) \triangleright (F', \phi') & \quad r' = \llbracket r / e \rrbracket_{\text{path}}^{\mathcal{E}, r} \\ \mathcal{E}; r'; s_1 \vdash \text{load } (G_1 ++ F' ++ G_2) \triangleright (v', d') \end{aligned}$$

By the induction hypothesis applied to  $s_1$  we have  $v' = v$  and  $\text{valid}(d) = \text{valid}(d')$ , which finishes the case.

**Case:**  $s = \langle x : s_1, s_2 \rangle$

By the definition of the load function we have

$$\begin{aligned} \mathcal{E}; r; s_1 \vdash \text{store } (F, v_1, d_1) \triangleright (F'_1, \phi'_1) & \quad v = (v_1, v_2) \\ (\mathcal{E}, x \mapsto v_1); r; s_2 \vdash \text{store } (F, v_2, d_2) \triangleright (F'_2, \phi'_2) & \quad d = (b, (d_1, d_2)) \\ & \quad F' = F'_1 ++ F'_2 \\ & \quad \phi' = \lambda F'. (b = \text{valid}(d_1) \wedge \text{valid}(d_2)) \wedge \phi'_1(F') \wedge \phi'_2(F') \end{aligned}$$

By  $\phi'(G_1 ++ F' ++ G_2)$  we have

$$\begin{aligned} b &= \text{valid}(d_1) \wedge \text{valid}(d_2) \\ \phi'_1(G_1 ++ F' ++ G_2) \\ \phi'_2(G_1 ++ F' ++ G_2) \end{aligned}$$

As  $(++)$  is associative we also have,

$$\begin{aligned} (G_1 ++ (F_1 ++ F_2) ++ G_2) &= (G_1 ++ F_1 ++ (F_2 ++ G_2)) \\ (G_1 ++ (F_1 ++ F_2) ++ G_2) &= ((G_1 ++ F_1) ++ F_2 ++ G_2), \end{aligned}$$

and hence:

$$\begin{aligned} \phi'_1(G_1 ++ F_1 ++ (F_2 ++ G_2)) \\ \phi'_2((G_1 ++ F_1) ++ F_2 ++ G_2) \end{aligned}$$

By the definition of the load function we have

$$\begin{aligned} \mathcal{E}; r; s_1 \vdash \text{load } (G_1 ++ F_1 ++ (F_2 ++ G_2)) \triangleright (v'_1, d'_1) & \quad v' = (v'_1, v'_2) \\ (\mathcal{E}, x \mapsto v_1); r; s_2 \vdash \text{load } ((G_1 ++ F_1) ++ F_2 ++ G_2) \triangleright (v'_2, d'_2) & \quad b' = \text{valid}(d_1) \wedge \text{valid}(d_2) \\ & \quad d' = (b', (d'_1, d'_2)) \end{aligned}$$



By the induction hypothesis applied to  $s_1$  and  $s_2$ , we have

$$\begin{aligned} v'_1 &= v_1 & \text{valid}(d_1) &= \text{valid}(d'_1) \\ v'_2 &= v_2 & \text{valid}(d_2) &= \text{valid}(d'_2) \end{aligned}$$

It follows that  $(v_1, v_2) = (v'_1, v'_2)$  and  $b = b'$ , which finishes the case.

**Case:**  $s = [s \mid x : e]$

By the definition of the store function we have

$$\begin{aligned} \forall i \in \{1, \dots, k\}. (\mathcal{E}, x \mapsto w_i); r; s \vdash \text{store } (F, v_i, d_i) \triangleright (F'_i, \phi_i) & & v &= [v_1, \dots, v_j] \\ \llbracket e \rrbracket_{[tau]}^{\mathcal{E}, r} = [w_1, \dots, w_m] & & d &= [d_1, \dots, d_l] \\ & & k &= \min(j, l, m) \\ & & F' &= F'_1 ++ \dots ++ F'_k \\ & & \phi &= \lambda F'. (j = k = l) \wedge (b = \bigwedge_i^k \text{valid}(d_i)) \bigwedge_i^k \phi_i(F') \end{aligned}$$

By  $\phi'(G_1 ++ F' ++ G_2)$  we have

$$\begin{aligned} k &= j = l = m & \forall i \in \{1, \dots, k\}. \phi'_i(G_1 ++ F' ++ G_2) \\ b &= (\bigwedge_i^k \text{valid}(d_i)) \end{aligned}$$

Let

$$H_i = ((G_1 ++ F_1 ++ \dots ++ F_{i-1}) ++ F_i ++ (F_{i+1} ++ \dots ++ F_k ++ G_2))$$

for  $i$  from 1 to  $k$ . As  $(++)$  is associative we have,

$$(G_1 ++ (F_1 ++ \dots ++ F_k) ++ G_2) = H_i \quad \text{for } i \in \{1, \dots, k\}$$

and hence:

$$\phi'_i(H_i) \quad \text{for } i \in \{1, \dots, k\}$$

By the definition of the load function we also have

$$\begin{aligned} \forall i \in \{1, \dots, k\} (\mathcal{E}, x \mapsto w_i); r; s_1 \vdash \text{load } H_i \triangleright (v_i, d'_i) & & v' &= [v'_1, \dots, v'_k] \\ & & b' &= \bigwedge_i^k \text{valid}(d'_i) \\ & & d' &= (b', [d'_1, \dots, d'_k]) \end{aligned}$$

By the induction hypothesis applied to  $s$  ( $k$  times), we have  $v'_i = v_i$  and  $\text{valid}(d'_i) = \text{valid}(d_i)$  for  $i$  from 1 to  $k$ . It follows that  $[v_1, \dots, v_k] = [v'_1, \dots, v'_k]$  and  $b = b'$ , which finishes the case.

**Case:**  $s = P(e)$

By the definitions of the store and load functions we have

$$\begin{aligned} v &= () & v' &= () \\ d &= ((), b) & d' &= (b', ()) \\ F' &= F & b' &= \llbracket e \rrbracket_{bool}^{\mathcal{E}, r} \\ \phi' &= \lambda F'. (b = \llbracket e \rrbracket_{bool}^{\mathcal{E}, r}) & & \end{aligned}$$

By  $\phi'(G_1 ++ F' ++ G_2)$  we have  $b = \llbracket e \rrbracket_{bool}^{\mathcal{E}, r}$ . It follows that  $v = v'$  and  $b = b'$ , which finishes the case.

**Case:**  $s = s_1?$

We analyze several subcases:

**Subcase:**  $v = \text{Just } v_1$  and  $d = (b, \text{Just } d_1)$

By the definition of the store function we have

$$\mathcal{E}; r; s_1 \vdash \text{store } (F, v_1, d_1) \triangleright (F', \phi'_1) \quad \phi' = \lambda F'. (b = \text{valid}(d_1)) \wedge (r \in \text{dom}(F)) \wedge \phi'_1(F')$$

By  $\phi'(G_1 ++ F' ++ G_2)$  we have  $b = \text{valid}(d_1)$  and  $r \in \text{dom}(G_1 ++ F' ++ G_2)$ . By the definition of the load function we also have

$$\begin{aligned} \mathcal{E}; r; s_1 \vdash \text{load } (G_1 ++ F' ++ G_2) \triangleright (v'_1, d'_1) & & v' &= \text{Just } v'_1 \\ & & b' &= \text{valid}(d'_1) \\ & & d' &= (b', \text{Just } d'_1) \end{aligned}$$

By the induction hypothesis applied to  $s_1$  we have  $v'_1 = v_1$  and  $b' = b$ . It follows that  $v' = v$ .

**Subcase:**  $v = \text{Nothing}$

By the definition of the store and load functions we have

$$\begin{aligned} F' &= F[r := \perp] & v' &= \text{Nothing} \\ \phi' &= \lambda F'. (d = \text{Nothing}) \wedge b \wedge (r \notin \text{dom}(F')) & b' &= \text{True} \\ & & d' &= (b', \text{Nothing}) \end{aligned}$$

As  $\phi'(G_1 ++ F' ++ G_2)$  we have  $d = \text{Nothing}$  and  $b$  and  $r \notin \text{dom}(G_1 ++ F' ++ G_2)$ . Thus, we immediately have  $v' = v$  and  $b' = b$ .

**Subcase:**  $v = \text{Just } v_1$  and  $d = \text{Nothing}$

Vacuously holds: by the definition of the store function we have  $\phi' = \lambda F'. \text{False}$ , which contradicts the assumption that  $\phi'(G_1 ++ F' ++ G_2)$ .

Thus, in each subcase we have  $v' = v$  and  $\text{valid}(d) = \text{valid}(d')$ , which finishes the subcase and the inductive proof. ■

## B. Appendix

This appendix contains contains Forest descriptions of a variety of different filestores. Please note that this appendix is best viewed electronically. Some of the graphs generated are very large, but shrunk down to fit on a single page. They will not display well when printed. However, reviewers may zoom in electronically on the PDF to view the details.

## C. Pads Web Site Description

This Forest description describes the Pads web site. The description starts with Pads descriptions of files that contain information that impacts the directory structure. The configuration file supplies the paths where various components of the website should be located. The SourceNames file lists the names of the data files available for the demo. Each user directory will have a subdirectory for each file listed in SourceNames. Each user is logged in the file UserEntries. For each user in this file, there is a directory with a corresponding name containing all of the information relevant to that user. A graph of the Pads website, generated using the ForestGraph tool follows the description.

```
[pads|
-- Configuration file for learning demo web site; contains paths to various web site components.
data Config_f = {
  header      :: [Pstringln] with term length of 13,
  "$host_name =", host_name  :: Config_entry_t,  --Name of machine hosting web site
  "$static_path =", static_path :: Config_entry_t, --URL prefix for static content
  "$cgi_path   =", cgi_path   :: Config_entry_t,  --URL prefix for cgi content
  "$script_path =", script_path :: Config_entry_t, --Path to directory of scripts in live web site
  "$tmp_root   =", tmp_root   :: Config_entry_t,  --Path to directory for demo user data
  "$pads_home  =", pads_home  :: Config_entry_t,  --Path to directory containing pads system
  "$learn_home =", learn_home :: Config_entry_t,  --Path to directory containing learning system
  "$sml_home   =", sml_home   :: Config_entry_t,  --Path to directory containing SML executable
  "$install_src =", install_src :: Config_entry_t, --Path to directory containing learning demo website source
  "$static_dst =", static_dst  :: Config_entry_t, --Path to directory for static content in live web site
  "$cgi_dst    =", cgi_dst    :: Config_entry_t,  --Path to directory for cgi content in live web site site
  trailer     :: [Pstringln]
}

type Config_entry_t = Line (" \"", Pstring '\\"', "\";")
type Header_t = [Pstringln] with term length of 13

{- File listing data sources for web site -}
type SourceNames_f = [Pstringln]

{- Information related to a single user's use of the web site -}
type UserEntries_f = [Line UserEntry_t] with term Eor

{- Each visitor gets assigned a userId that is passed as a ? parameter in URL.
Security considerations preclude using user-modifiable values as part of file paths.
Thus, we map each userId to a corresponding dirId.
The dirId names the directory containing the associated user's data.
A userEntry_t contains a single such mapping.
A file with type userEntries_t describes a collection of such mappings.
-}
data UserEntry_t = {
  "id.",   usrId :: Pint,
  ",id.",  dirId :: (Pint, '.', Pint)   where <| usrId == fst dirId |>
}

{- Log of requests. Used to prevent denial of service attacks. -}
type LogFile_f = [LogEntry_t]
```

```

{- Request entry. -}
data LogEntry_t = {
  userId :: Pint,           '','',' --user making request
  ip      :: IP_t,         '','',' --IP address of requestor
  script  :: Pstring '','',' --script to be executed
  userDir :: Pstring '','',' --directory to put results, corresponds to user
  padsv   :: Pstring '','',' --version of PADS used
  sml     :: PstringSE(RE " "), --version of SML used
  msg     :: Maybe Pstringln --optional message
}

type IP_t = (Pint, '.', Pint, '.', Pint, '.', Pint)
[]

[forest]
{- Files with various permission settings. -}
type BinaryRO = Binary      where <| get_modes this_att == "-rw-r--r--" |>
type BinaryRX = Binary      where <| get_modes this_att == "-rwxr-xr-x" |>
type TextRX   = Text        where <| get_modes this_att == "-rwxr-xr-x" |>
type TextRO   = Text        where <| get_modes this_att == "-rw-r--r--" |>

{- Optional binary file with read/execute permission. -}
type OptBinaryRX = Maybe BinaryRX

{- Files with PADS descriptions -}
type Config = File Config_f      where <| get_modes this_att == "-rw-r--r--" |>
type SourceNames = File SourceNames_f where <| isReadOnly this_att |>
type UserEntries = File UserEntries_f where <| isReadOnly this_att |>
type LogFile = File LogFile_f    where <| isReadOnly this_att |>

{- Directory of image files -}
type Imgs_d = Directory {
  logo    is "pads_small.jpg" :: BinaryRO,
  favicon is "favicon.ico"   :: BinaryRO
}

{- Directory of static content -}
type Static_d = Directory {
  style_sheet is "pads.css"           :: TextRO,
  intro_redir is "learning-demo.html" :: TextRO,
  title_frame is "atitle.html"       :: TextRO,
  logo_frame  is "top-left.html"     :: TextRO,
  top_frame   is "banner.html"       :: TextRO,
  empty_frame is "nothing.html"      :: TextRO,
  images      is "images"            :: Imgs_d where <| get_modes images_md == "drwxr-xr-x" |>
}

{- Directory of dynamic content -}
type Cgi_d = Directory {
  config' is "PLConfig.pm"           :: TextRO,
  perl_utils is "PLUtilities.pm"     :: TextRO,
  intro     is "learning-demo.cgi"   :: TextRX,
  intro_nav is "navbar-orig.cgi"     :: TextRX,
  select_data is "pads.cgi"          :: TextRX,
  result_nav is "navbar.cgi"         :: TextRX,
  format_chosen is "data-results.cgi" :: TextRX,
  gen_desc    is "build-description.cgi" :: TextRX,
  get_user_data is "build-roll-your-own.cgi" :: TextRX,
  gen_desc_usr is "genData.cgi"      :: TextRX,
  build_lib   is "build-library.cgi"  :: TextRX,
  build_accum is "build-accum.cgi"    :: TextRX,
  build_xml   is "build-xml.cgi"     :: TextRX,
  build_fmt   is "build-fmt.cgi"     :: TextRX
}

```

```

{- Directory of shell scripts invoked by CGI to run learning system -}
type Scripts_d = Directory {
  rlearn      :: TextRX,           --Shell script for running PADS comiler on stock format
  rlearnown  is "rlearn-own" :: TextRX, --Shell script for running PADS compiler on user format
  raccum    is "r-accum"   :: TextRX, --Shell script to generate and run accumulator
  rxml       is "r-xml"    :: TextRX, --Shell script to generate and run XML converter
  rfmt       is "r-fmt"    :: TextRX, --Shell script to generate and run formatting program
  rlibrary   :: TextRX           --Shell script to build PADS library
}

{- Directory containing administrative files used by demo web site -}
type Info_d = Directory {
  sources is "sampleFiles" :: SourceNames, --List of source data files whose formats can be learned
  users   is "userFile"    :: UserEntries, --Mapping from userIDs to associated directory names
  logfile is "logfile"     :: LogFile      --Log of server actions
}

{- Collection of files named by sources containing actual data. -}
type DataSource_d(sources :: [String]) = [ s :: Text | s <- sources ]

{- Type of a symbolic link with pointing to source-}
type SymLink_f (path :: FilePath) = SymLink where <| this == path |>

{- Directory of optional links to source data files -}
type Data_d ((root,sources) :: (FilePath, [String])) = Directory {
  datareps is [s :: Maybe Text | s <- sources],
  datalinks is [s :: Maybe (SymLink_f <| root++"/"++ s |>) | s <- sources]
}

{- Directory that stores the generated machine-dependent output for data source named source -}
type MachineDep_d (source :: String) = Directory {
  pads_c   is <| source ++ ".c"   |> :: TextRO,           --Generated C source for PADS description
  pads_h   is <| source ++ ".h"   |> :: TextRO,           --Generated C header for PADS description
  pads_o   is <| source ++ ".o"   |> :: BinaryRO,         --Compiled library for PADS description
  pads_pxml is <| source ++ ".pxml"|> :: TextRO,           --PADS description in xml syntax
  pads_xsd is <| source ++ ".xsd" |> :: TextRO,           --Xschema of XML syntax for source description
  pads_acc is <| source ++ "-accum"|> :: OptBinaryRX,     --Optional generated accumulator program
  pads_fmt is <| source ++ "-fmt" |> :: OptBinaryRX,     --Optional generated formatting program
  pads_xml is <| source ++ "-xml" |> :: OptBinaryRX      --Optional generated XML conversion program
}

{- Directory that stores the generated output for data source named "source". -}
type Example_d (source :: String) = Directory {
  pads_p      is <| source ++ ".p" |> :: TextRO,           --PADS/C description of data source
  pads_pml    is <| source ++ ".pml" |> :: Maybe TextRO,  --PADS/ML description of data source
  vanilla     is "vanilla.p"      :: TextRO,           --input tokenization
  makefile    is "GNUmakefile"    :: Text,            --Makefile
  machine     is <| envVar "AST_ARCH"|> :: Maybe (MachineDep_d source), --Platform dependent files
  accum_c     is <| source ++ "-accum.c" |> :: Maybe TextRO, --Template for accumulator program
  accum_out   is <| source ++ "-accum.out"|> :: Maybe TextRO, --ASCII Accumulator output
  accum_xml_out is <| source ++ "-accum_xml.out"|> :: Maybe TextRO, --XML Accumulator output
  xml_c       is <| source ++ "-xml.c"|> :: Maybe TextRO,  --Template for XML converter
  xml_out     is <| source ++ "-xml.out"|> :: Maybe TextRO, --XML representation of source
  xml_xsd     is <| source ++ ".xsd" |> :: Maybe TextRO,  --Xschema for XML representation of source
  fmt_c       is <| source ++ "-fmt.c" |> :: Maybe TextRO, --Template for formatting program
  fmt_out     is <| source ++ "-fmt.out" |> :: Maybe TextRO --Formatted representation of source
}

{- Directory that stores all information for one user. -}
type User_d(arg@ (r, sources) :: (FilePath, [String])) = Directory {
  dataSets is "data" :: Maybe (Data_d arg),
  runExamples is [ s :: Maybe (Example_d s) | s <- sources]
}

```

```

{- Collection of directories containing temporary information for all users. -}
type Users_d((r,info) :: (FilePath, Info_d)) =
  [userDir :: User_d <|(r, getSources info) |> | userDir <- <| userNames info |> ]

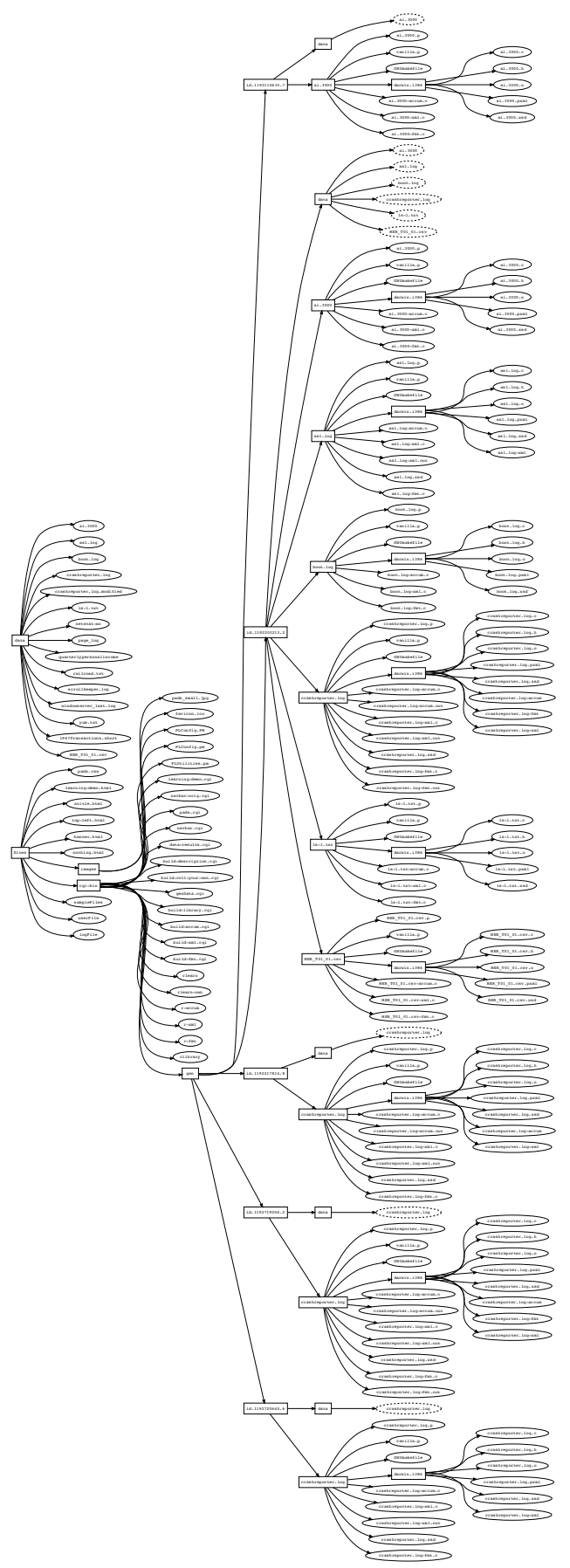
{- Top-level of PADS website. -}
type Website_d(config::FilePath) = Directory {
  c
  static_content is <| gstatic_dst c |> :: Config,           --Configuration file with locations of other components
  dynamic_content is <| gcgi_path c |> :: Static_d,           --Static web site content
  scripts         is <| gscript_path c |> :: Cgi_d,           --Dynamic web site content
  admin_info      is <| gstatic_dst c |> :: Scripts_d,        --Shell scripts invoked by cgi to run learning system
  data_dir        is <| (glearn_home c)++"/examples/data" |>  --Administrative information about website
  usr_data        is <| (gtmp_root c)++"/examples/data" |>    --Stock data files for website
  }
  : DataSource_d <|(getSources admin_info)|>,
  : Users_d <|(get_fullpath data_dir_md, admin_info)|>  --User info
}
[]

{- HASKELL HELPER FUNCTIONS -}
isReadOnly md = get_modes md == "-rw-r--r--"

{- Function userName gets the list of user directorn names from an info structure. -}
userNames info = getUserEntries (users info)
getUserEntries (UserEntries (UserEntries_f users)) = map userEntryToFileName users
userEntryToFileName userEntry = pairToFileName (dirId userEntry)
pairToFileName (Pint n1, Pint n2) = "id."++(show n1)++"."++(show n2)

{- Helper functions to convert a Config entry to a FilePath -}
cToS (Config_entry_t (Pstring s)) = s
ghost_name (Config c) = cToS $ host_name c
gstatic_path (Config c) = cToS $ static_path c
gcgi_path (Config c) = cToS $ cgi_path c
gscript_path (Config c) = cToS $ script_path c
glearn_home (Config c) = cToS $ learn_home c
gtmp_root (Config c) = cToS $ tmp_root c
gstatic_dst (Config c) = cToS $ static_dst c

```



## D. Students.hs Description

This section includes the Forest description of the Princeton Computer Science Department filestore. The following is the initial portion of a student record, shown here to illustrate the format.

```
KESSEL, PHIL          BSE '11
-----
Type   Yr  Course      Grade
      1          A+ to F
d      2          P ( Pass )
t D p  3          INC
o . .  4 Dept xxx N (Not Avail)
-----
d . .  1 COS   101 C
o . .  1 HOC   101 A
o . .  1 GOL   599 A+
...
```

```
-- Auxiliary Haskell functions for PADS description
```

```
ws  = RE "[ ]+"
ows = RE "[ ]*"
junk = RE ".*"
space = ' '
quote = '"'
comma = ','
```

```
-- PADS description of Princeton CS Student Record Format
```

```
[pads|
  type Grade = Pre "[ABCD][+-]?|F|AUD|N|INC|P"

  data Course =
    { sort           :: Pre "[dto]",           ws
    , departmental  :: Pre "[.D]",           ws
    , passfail      :: Pre "[.p]",           ws
    , level         :: Pre "[1234]",         ws
    , department    :: Pre "[A-Z][A-Z][A-Z]", ws
    , number        :: Pint where <| 100 <= number && number < 600 |>, ws
    , grade         :: Grade,                junk
    }

  data Middle_name = {space, middle :: Pre "[a-zA-Z][.]" }

  data Student_Name(myname::String) =
    { lastname      :: Pre "[a-zA-Z]*" where <| toString lastname == myname |>, comma, ows
    , firstname     :: Pre "[a-zA-Z]*"
    , middlename    :: Maybe Middle_name
    }

  data School = AB | BSE

  data Person (myname::String) =
    { fullname      :: Student_Name myname,    ws
    , school        :: School,                 ws, quote
    , year          :: Pre "[0-9][0-9]"
    }

  type Header = [Line (Pre ".*")] with term length of 7
  type Trailer = [Line (Pre ".*")] with term Eof
  data Student (name::String) =
    { person       :: Line (Person name)
    , Header
    , courses      :: [Line Course]
    , Trailer
    }
|]
```

```

-- Auxiliary Haskell functions for Forest description
template s = s `elem` ["SSSS.txt", "SSS.txt", "sxx.txt", "sss.txt", "ssss.txt"]
not_template = not . template

getYear :: String -> Integer
getYear s = read (reverse (take 2 (reverse s)))
toStrN i n = (replicate (n - length (show i)) '0') ++ (show i)
mkClass y = "classof" ++ (toStrN y 2)

transferRE = RE "TRANSFER|Transfer"
leaveRE    = RE "LEAVE|Leave"
withdrawnRE = RE "WITHDRAWN|WITHDRAWAL|Withdrawn|Withdrawal|WITHDREW"
cRE        = RE "classof[0-9][0-9]"
txt         = GL "*.txt"

-- FOREST description of Princeton CS Department Database
[forest]
-- Root of the hierarchy
type PrincetonCS (y::Integer) = Directory
  { notes is "README" :: Text
  , seniors is <|mkClass y |> :: Class y
  , juniors is <|mkClass (y + 1)|> :: Class <| y + 1 |>
  , graduates :: Grads
  }

-- Collection of directories containing graduated students
type Grads =
  Map [ c :: Class <| getYear c |> | c <- matches cRE ]

-- Directory containing all students in a particular year
type Class (y :: Integer) = Directory
  { bse is <|"BSE" ++ (toStrN y 2)|> :: Major
  , ab is <|"AB" ++ (toStrN y 2)|> :: Major
  , transfer matches transferRE :: Maybe Major
  , withdrawn matches withdrawnRE :: Maybe Major
  , leave matches leaveRE :: Maybe Major
  }

-- Collection of files containing all students in a particular major.
type Major = Map
  [ s :: File (Student <| dropExtension s |>)
  | s <- matches txt, <| (not . template) s |> ]
[]

```



## D.1 Generated Description

Here follows a description generated from a small sample of the student directory data using the description inference tool.

```
data transfer = Directory {
}
data WITHDREW = Directory {
  fingertxt is "finger.txt" :: File Ptext
}
data tTRANSFER = Directory {
  BEAUCHEMINtxt is "BEAUCHEMIN.txt" :: File Ptext,
  VERSTEEGtxt is "VERSTEEG.txt" :: File Ptext
}
data bSE11 = Directory {
  transfer is "transfer" :: transfer,
  BOZAKtxt is "BOZAK.txt" :: File Ptext,
  KESSELTtxt is "KESSEL.txt" :: File Ptext,
  ssstxt is "sss.txt" :: File Ptext
}
data aB11 = Directory {
  KADRItxt is "KADRI.txt" :: File Ptext,
  MACARTHERtxt is "MACARTHER.txt" :: File Ptext,
  ORRtxt is "ORR.txt" :: File Ptext,
  SSSStxt is "SSSS.txt" :: File Ptext
}
data classof11 = Directory {
  aB11 is "AB11" :: aB11,
  bSE11 is "BSE11" :: bSE11,
  tTRANSFER is "TRANSFER" :: tTRANSFER,
  WITHDREW is "WITHDREW" :: WITHDREW
}
```

## E. Coral.hs Description

This section gives the PADS and Forest descriptions for the CoralCDN Log repository. A graph of the CoralCDN repository, generated like the graph above using the ForestGraph tool from the description and (a subset of) the actual repository follows.

```
-- Auxiliary Haskell definitions for PADS description
comma_ws = RE ", [ ]*"
status_re = RE "[0-9]+"

-- PADS description of CoralCDN Webserver Log Format
[pads|
type Time = (Pint, ".", Pint)

type Byte = constrain x :: Pint where <| 0 <= x && x <= 256 |>

type IP_Port =
  { "",
    ip :: (Byte, '.', Byte, '.', Byte, '.', Byte), ":",
    port :: Pint, "" }

type Status = PstringME(status_re)

type Statistics =
  { stats_size      :: Pint,      comma_ws
  , stats_proxy    :: Pre "[01]", comma_ws
  , stats_level    :: Pint,      comma_ws
  , stats_lookup   :: Pint,      comma_ws
  , stats_xfer     :: Pint,      comma_ws
  , stats_total    :: Pint }

type NoQuote = PstringME (RE "[^\\"]*")

type Generic = ('"', NoQuote, "'")

type Url = Generic

data Header =
  { version      :: Maybe (Pre "[12],[ \\t]*")
  , time        :: Time      }

data Request =
  { src          :: IP_Port, comma_ws
  , dst         :: IP_Port, comma_ws
  , url         :: Url    }

data InData =
  { "\\IN\\",      comma_ws
  , in_req       :: Request, comma_ws
  , in_status1  :: Status,  comma_ws
  , in_status2  :: Status,  comma_ws
  , in_stats    :: Statistics }

data OutData =
  { "\\OUT\\",      comma_ws
  , out_remote  :: Pre "\\(REM|LOC)\\", comma_ws
  , out_req     :: Request,  comma_ws
  , out_referrer :: Url,     comma_ws
  , out_status  :: Status,   comma_ws
  , out_stats   :: Statistics, comma_ws
  , out_forwarded :: Generic, comma_ws
  , out_via     :: Generic   }

data InOut = In InData | Out OutData

data Entry =
  { header :: Header,  comma_ws
  , payload :: InOut
  , Eor }

type Entries = [Entry] with term Eor

type Coral = (Entries, Eof)
|]
```

```
-- Forest description of CoralCDN Log Repository
[forest]
-- Directory containing log files
type Log = Directory
  { web is "coralwebsrv.log.gz" :: Gzip (File Coral),
    dns is "coraldnssrv.log.gz" :: Maybe (Gzip (File Ptext)),
    prb is "probed.log.gz"      :: Maybe (Gzip (File Ptext)),
    dmn is "corald.log.gz"     :: Maybe (Gzip (File Ptext)) }

-- Directory containing dates
type Site = [ d :: Log | d <- matches (RE "[0-9]{4}_[0-9]{2}_[0-9]{2}-[0-9]{2}_[0-9]{2}") ]

-- Directory containing sites
type Top = [ s :: Site | s <- matches (RE "[^\\.]*") ]
[]
```

```

-- Load function for CoralCDN description
(rep,md) = unsafePerformIO $ top_load "/var/log/coral"

-- Helpers: deconstruct representations
get_sites :: Top -> [(String,Site)]
get_dates :: Site -> [(String,Log)]
get_entries :: Log -> [Entry]

-- Helpers: project fields
get_stats :: Entry -> Statistics
get_total :: Entry -> Int
get_date :: String -> String
get_url::Entry -> String
string_of_url :: Url -> String
is_in :: Entry -> Bool
is_out :: Entry -> Bool

-- Helper: builds an association list
lmap f p tdir =
  [ f host datetime e | (host,hdir) <- get_sites tdir,
                        (datetime,ldir) <- get_dates hdir,
                        e <- get_entries ldir,
                        p e ]

-- Uses of lmap
by_date = lmap (\h d e -> (get_date d, get_total e))
by_host = lmap (\h d e -> (h, get_total e))
by_url_bytes = lmap (\h d e -> (get_url e, get_total e))
by_url_counts = lmap (\h d e -> (get_url e, 1))

-- Helpers: fold down an association list
go_bins m p = fromListWith (+) (m p rep)

count_bins m = fromListWith (+) (fold (\ c l -> (c,l):l) [] m)

go_flat p =
  sum [ (get_total e) | (host,hdir) <- get_sites tdir,
                       (datetime,ldir) <- get_dates hdir,
                       e <- get_entries ldir,
                       p e ]

-- Several useful queries
in_total = go_flat is_in
out_total = go_flat is_out
in_by_host = go_bins by_host is_in
out_by_host = go_bins by_host is_out
in_by_date = go_bins by_date is_in
out_by_date = go_bins by_date is_out
in_url_bytes = go_bins by_url_bytes is_in
out_url_bytes = go_bins by_url_bytes is_out
in_url_counts = go_bins by_url_counts is_in
out_url_counts = go_bins by_url_counts is_out
in_counts_urls = count_bins $ go_bins by_url_counts is_in
out_counts_urls = count_bins $ go_bins by_url_counts is_out
num_sites () = case load_logs () of Top l -> List.length l

-- Top-k URLs
topk k =
  take k $ sortBy sortDown $ toList $
  fromListWith (+)
  [ (get_url e, get_total e)
  | (site,sdir) <- get_sites rep,
    (datetime,ldir) <- get_dates sdir,
    e <- get_entries ldir,
    is_in e ]

```



## F. Gene Ontology

This section presents a description of gene ontology data found here: <http://www.geneontology.org/gene-associations/>. A graph generated using ForestGraph on a subset of the data follows the description.

This filestore is a web directory of gene association data files. The root directory contains a number of .gz files, a readme directory and a submission directory. Each .gz file is the gene ontology (GO) data of the genes in one or more organism, and the file names have the format "gene\_association.XXX\_YYY.gz", where XXX represents the name of the institute that provides the data and YYY is the name of the organism. YYY is optional because some institute provides the data for only one organism.

The readme directory contains a set of .README files for a subset of the GO data in the root.

The submission directory contains a set of .gz files, their corresponding .conf files, and a paint sub-directory. The .gz files are similar to the ones in root except they are older. The .conf file summarizes some attributes of the .gz file such as "the name of the project", "contact email", etc. The paint sub-directory contains a further set of subdirectories of the form PTHRXXXXXX, where XXXXX is a 5-digit number. These subdirectories each contain six text files and an XML file. These are the annotation inference of the gene ontology using phylogenetic trees and the PAINT tool.

```
-- PADS descriptions of data file format.
[ pads |
  type Pfloat          = (Pint, '.', Pint)
  type Pdate           = {mon :: Pint, '/', day :: Pint, '/', year :: Pint}
  type Purl            = ("http://", Pstringln)
  type Version_t       = ("!CVS Version: Revision: ", Pfloat, ws, '$')
  type Valid_date_t    = ("!GOC Validation Date: ", Pdate, ws, '$')
  type Sub_date_t      = ("!Submission Date: ", Pdate)
  type Project_name_t  = ("!Project_name: ", Pstringln)
  type URL_t           = ("!URL: ", Purl)
  type Email_t         = ("!Contact Email: ", Pstringln)
  type Funding_t       = ("!Funding: ", Pstringln)
  type Gaf_ver_t       = ("!gaf-version: ", Pfloat)
  type Organism_t      = ("!organism:", ws, Pstringln)
  type Date_t          = ("date:", ws, Pdate)
  type Note_t          = ('!', ws, Pstringln)

  data Header_line_t =
    Version Version_t
    | Valid_date Valid_date_t
    | Sub_date Sub_date_t
    | Project_name Project_name_t
    | URL URL_t
    | Email Email_t
    | Funding Funding_t
    | Gaf_ver Gaf_ver_t
    | Organism Organism_t
    | Date Date_t
    | Note Note_t
    | Other ('!', Pstringln)

  type Other_line_t = Pstringln

  type GA_f = ([Line Header_line_t], [Line Other_line_t] with term Eof)
]

[ pads |
  data Pair_t = {key::Pstring '=', '=', val::Pstringln}
  type Conf_f = [Line Pair_t] with term Eof
]

[ pads |
  type Xml_header = ("<?xml ", Pstringln)
  type XML_f = (Line Xml_header, [Line Pstringln])
]
```

```

-- Forest description of Gene Ontology filestore
[forest]
  type Readme_d = Directory {
    readmes is [rm :: Maybe Text | rm <- <|map get_readme_file (comb_source sources)|>]
  }

  type PTHR_d (name :: String) = Directory {
    attr is <| name ++ ".save.attr" |> :: Text,
    gaf is <| name ++ ".save.gaf" |> :: Text,
    msa is <| name ++ ".save.msa" |> :: Text,
    paint is <| name ++ ".save.paint" |> :: File XML_f,
    sfan is <| name ++ ".save.sfan" |> :: Text,
    tree is <| name ++ ".save.tree" |> :: Text,
    txt is <| name ++ ".save.txt" |> :: Text,
    wts is <| name ++ ".save.txt" |> :: Text
  }

  type Pre_sub_d = Directory {
    pre_gz_files is [gz :: Maybe (Gzip (File GA_f)) | gz <- <|map get_gz_file (comb_source sources)|>],
    pre_conf_files is [conf :: Maybe (File Conf_f) | conf <- <|map get_conf_file (comb_source sources)|>]
  }

  type Paint_d = Directory {
    pthr_dirs is [dir_name :: PTHR_d (dir_name) | dir_name <- matches RE "PTHR[0-9]+"],
    pre_sub is "pre-submission" :: Pre_sub_d
  }

  type Submission_d = Directory {
    gz_files is [gz :: Maybe (Gzip (File GA_f)) | gz <- <|map get_gz_file (comb_source sources)|>],
    conf_files is [conf :: Maybe (File Conf_f) | conf <- <|map get_conf_file (comb_source sources)|>],
    paint_files is [cs :: Maybe (File Conf_f)
                   | cs <- <|map (\x -> get_conf_file ("paint" ++ x)) (comb_source sources)|>],
    paint_d is "paint" :: Paint_d
  }

  type Top_d = Directory {
    data_files is [gz :: Maybe (Gzip (File GA_f)) |
                  gz <- <|map get_gz_file (comb_source sources)|>],
    readme is "readme" :: Readme_d,
    sub is "submission" :: Submission_d
  }
]

-- Haskell code to generate graph corresponding to sample data set in filestore "Data/ga"
doImg = do
  (rep,md) <- top_d_load "Data/ga"
  ; mdToPDF md "Examples/ga.pdf"

```

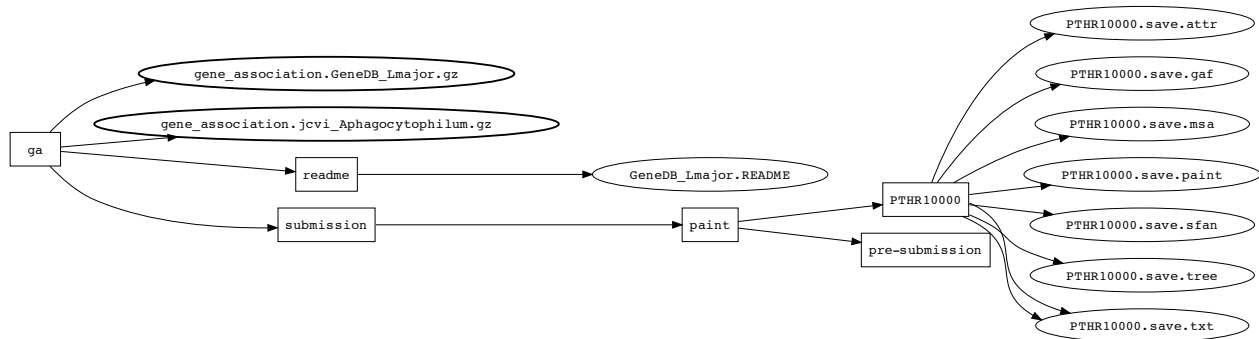
```

-- Auxiliary Haskell Definitions
ws = RE "[ ]+"
title = "gene_association"
get_gz_file f = title ++ "." ++ f ++ ".gz"
get_readme_file f = f ++ ".README"
get_conf_file f = title ++ "." ++ f ++ ".conf"

{- each source is a pair (institute name, list of organisms the institute provides) -}
sources = [
  ("Compugen", [])
  , ("GeneDB", ["Lmajor", "Pfalciparum", "Spombe", "Tbrucei", "tsetse"])
  , ("PAMGO", ["Atumefaciens", "Ddadantii", "Mgrisea", "Oomycetes"])
  , ("aspgd", [])
  , ("cgd", [])
  , ("dictyBase", [])
  , ("ecocyc", [])
  , ("fb", [])
  , ("goa", ["arabidopsis", "chicken", "cow", "human", "mouse", "pdb", "rat",
            "uniprot", "uniprot_noiea", "zebrafish"])
  , ("gramene", ["oryza"])
  , ("jcbi", ["Aphagocytophilum", "Banthraxis", "Cburnetii", "Chydrogenoformans",
            "Cjejuni", "Cperfringens", "Cpsychrerythraea", "Dethenogenes", "Echaffeensis",
            "Gsulfurreducens", "Hneptunium", "Lmonocytogenes", "Mcapsulatus", "Nsennetsu",
            "Pfluorescens", "Psyringae", "phaseolicola", "Soneidensis", "Spomeroyi",
            "Vcholerae"])
  , ("mgi", [])
  , ("pseudocap", [])
  , ("reactome", [])
  , ("rgd", [])
  , ("sgd", [])
  , ("sgn", [])
  , ("tair", [])
  , ("wb", [])
  , ("zfin", []) ]
comb_source [] = []
comb_source ((inst, organs):sources) =
  let cl = case organs of
        [] -> [inst]
        _ -> map (\organism -> inst ++ "_" ++ organism) organs
  in cl ++ (comb_source sources)

{- the GO files, when unzipped, contain a header like the following:
!CVS Version: Revision: 1.19 $
!GOC Validation Date: 01/27/2007 $
!Submission Date: 1/15/2007
-}

```





## G. CVS.hs Description

This section provides a generic description for CVS repositories.

```
-- PADS description of CVS file formats
[ pads |
  type Repository_f = Line Pstringln
  data Mode_t = Ext ":ext:" | Local ":local:" | Server ":server:"
  data Root_t = { cvs_mode :: Maybe Mode_t
                , machine  :: Pstring ':', ':'
                , path     :: Pstringln
                }
  type Root_f = Line Root_t
  data Dentry_t = { "D/"
                  , dirname :: Pstring '/'
                  , "////"
                  }

  data Revision_t = Version (Pint, '.', Pint) | Added '0' | Removed '-'
  data TimeStamp_t = { ts          :: PstringSE (RE "[/+]")
                     , conflict  :: Maybe ('+', Pstring '/') }

  type Fentry_t = {
                  , filename  :: Pstring '/',
                  , revision  :: Revision_t,
                  , timestamp :: TimeStamp_t,
                  , options   :: Pstring '/',
                  , tagdate   :: Pstringln
                  }

  data Entry_t = Dir Dentry_t | File Fentry_t | NoDir 'D'
  type Entries_f = [Line Entry_t] with term Eof
]

-- Auxiliary Haskell functions
getEntries cvs = let (Entries_f l) = entries cvs in l
getDirName d = let (Pstring s) = dirname d in s
getFileName f = let (Pstring s) = filename f in s

isDir entry = case entry of Dir _ -> True; otherwise -> False
isFile entry = case entry of File _ -> True; otherwise -> False

getDirs cvs = map (\(Dir d) -> d) (filter isDir (getEntries cvs))
getFiles cvs = map (\(File f) -> f) (filter isFile (getEntries cvs))

-- FOREST description of CVS directory structure
-- Note that this description is recursive.
-- Note also that the collection of dirs and the
-- collection of files are determined from information in the cvs
-- directory.
[ forest |
  type CVS_d = Directory
  { repository is "Repository" :: File Repository_f
  , root       is "Root"       :: File Root_f
  , entries    is "Entries"    :: File Entries_f
  }

  type CVS_Repository_d = Directory
  { cvs          is "CVS"          :: CVS_d
  , dirs         is [ n as <| getDirName d |> :: CVS_Repository_d | d <- <| getDirs cvs |> ]
  , files        is [ <| getFileName f |> :: Text          | f <- <| getFiles cvs |> ]
  } []

-- Sample use of PADS and FOREST descriptions
meta_dir = "Examples/ CVS "
entries_file = meta_dir ++ "/Entries"
doParseEntries = do {
  (rep, md) <- parseFile entries_file
}

doLoadCVS = do {
  (meta_rep, meta_md) <- cvs_d_load meta_dir
}

```

## H. Universal.hs Description

This section includes a universal data description. This universal description is used to drive some of our generic tools.

```
-- Universal Forest Directory Description
```

```
[forest|
  type Universal_d = Directory
  { ascii_files  is [ f :: Text          | f <- matches (GL "*"), <| get_kind f_att == AsciiK    |> ]
  , binary_files is [ b :: Binary       | b <- matches (GL "*"), <| get_kind b_att == BinaryK   |> ]
  , directories  is [ d :: Universal_d | d <- matches (GL "*"), <| get_kind d_att == DirectoryK |> ]
  , symLinks     is [ s :: SymLink      | s <- matches (GL "*"), <| get_isSym s_att == True    |> ]
  }
|]
```

```
-- Use of Universal directory
```

```
universal_dir = "Examples/data/universal"
doLoadUniverse = do {
  (rep, md) <- universal_d_load universal_dir
}
```