

9 10 10.95 12

A Specification for the Synergy File System

William R. Bevier
Richard Cohen

Jeff Turner ¹

Technical Report 120

September, 1995

Computational Logic, Inc.
1717 West Sixth Street, Suite 290
Austin, Texas 78703-4776

TEL: +1 512 322 9951
FAX: +1 512 322 0656

EMAIL: bevier@cli.com, cohen@cli.com, sjt@tycho.ncsc.mil

Copyright © 2004 Computational Logic, Inc.

¹National Security Agency, R23

Contents

1	Introduction	1
2	Preliminaries	2
3	The Basic File Space	4
4	The File Space	6
5	The File Table	14
6	Processes	19
7	Garbage Collection	21
8	The Basic Access System	22
9	Discretionary Access Control	27
10	The Unix File System Interface	30
11	Mandatory Access Controls	36
12	The Synergy File System Interface	38
13	Reactions to the Z notation	42
14	Errors Uncovered by ACL2 Model	43

1 Introduction

This document contains a Z specification [Spivey, 1989] for a subset of the user interface functions of the Synergy file system. The Synergy file system is intended to be compatible with the Unix file system, but include additional constraints on behavior related to security. Morgan and Sufrin have previously specified the Unix file system interface [Morgan and Sufrin, 1987]. We have re-specified this interface to include features they did not treat, for example, discretionary access control. Our specification is influenced by the descriptions of the Unix file system given by *The Design of the Unix Operating System* [Bach, 1986] and *The Design and Implementation of the 4.3BSD UNIX Operation System* [Leffler et al., 1989]. Additional details have been gleaned from the *UNIX User's Reference Manual* for 4.3BSD [CSRG, 1986] and Steven's treatise on Unix programming [Stevens, 1992].

The work described by this paper was performed as a subproject of the Synergy project, whose primary goal is to develop policy-flexible security architectures that meet the needs of both the DoD (separation of data with high assurance) and the commercial world (flexibility, larger market than DoD). This work will serve as a foundation for several future projects. First, we will continue to articulate this specification in order to fully capture the design and interface of a policy-independent file system interface. Ultimately, this document should server as a programmer's manual for the file system. Second, we will use it as the highest level specification in a proof that the decomposition of the interface into a number of independent servers is valid.

In building a model of a system one must make choices regarding level of abstraction. These choices reveal properties of interest and ignore others. We are interested in mandatory and discretionary access controls on files. We are less interested in covert channels that may be created by these mechanisms. We therefore ignore limitations on virtually all resources. Such limitations can be introduced gradually to allow analysis of the channels they produce.

We have left error reporting and return codes out of our model. The model does specify the preconditions necessary for successful operations within the file system. We simply do not describe the effect of an operation if its preconditions are not met.

Section 2 of this report contains definitions of some primitive Z functions that are used elsewhere. Subsequent sections incrementally specify aspects of the file system interface. The following table summarizes the concepts that are introduced.

State	Concepts Introduced	Operations Introduced
Basic File Space	files, file identifiers	create, destroy, read, write
File Space	directories, pathnames	link, unlink
File Table	access modes, file position	open, close
Processes	processes, file descriptors	fork process, duplicate fd
Garbage Collection	reclaiming inaccessible files	
Basic Access System		
DAC	user, group, file permissions	
Unix FS Interface	current process, working directory DAC checks for operations	
MAC	security context	
Synergy FS Interface	MAC checks for operations	get/set context

2 Preliminaries

We presume the reader is familiar with the Z notation (e.g., has access to the *Z Reference Manual* [Spivey, 1989]).

Naming Conventions

In general the words “create” and “destroy”, “link” and “unlink” are used to describe file operations. The words “add” and “delete” are used to describe operations on other data structures/state components. One exception is the use of “create” and “destroy” to name the operations that create and destroy file and process attributes in the DAC and MAC layers.

Some Z preliminaries

In the Z notation, a sequence of length \mathbf{n} is a mapping from natural numbers between 1 and \mathbf{n} to the individual elements of the sequence. That is, sequence indexing is 1-origin.

The function **Nthtail** returns the subsequence of a sequence beginning at a given offset. Essentially, it skips \mathbf{n} things, and gives a new sequence starting there. Note that for all sequences \mathbf{s} , $\mathbf{nthtail}(0, \mathbf{s}) = \mathbf{s}$. That is, an offset of 0 skips nothing at the beginning of a sequence. This is equivalent to Morgan and Sufrin’s function **after**.

$\mathbf{[X]}$
nthtail : $\times \text{seq } \mathbf{X} \rightarrow \text{seq } \mathbf{X}$
$\forall \mathbf{n} ; ; \mathbf{s} : \text{seq } \mathbf{X} \bullet$ $\text{dom}(\mathbf{nthtail}(\mathbf{n}, \mathbf{s})) = (1 \dots \#\mathbf{s} - \mathbf{n}) \wedge$ $(\forall \mathbf{i} : \bullet$ $\quad (\mathbf{i} + \mathbf{n}) \in \text{dom } \mathbf{s} \Rightarrow (\mathbf{nthtail}(\mathbf{n}, \mathbf{s}))(\mathbf{i}) = \mathbf{s}(\mathbf{i} + \mathbf{n}))$

Shift adjusts the domain of a sequence forward by some offset. Morgan and Sufrin define the same function, but make it an infix operator.

$\mathbf{[X]}$
shift : $\times \text{seq } \mathbf{X} \rightarrow (\mathbf{X})$
$\forall \mathbf{n} ; ; \mathbf{s} : \text{seq } \mathbf{X} \bullet$ $\text{dom}(\mathbf{shift}(\mathbf{n}, \mathbf{s})) = \{\mathbf{i} : \text{dom } \mathbf{s} \bullet \mathbf{i} + \mathbf{n}\} \wedge$ $(\forall \mathbf{i} : \text{dom}(\mathbf{shift}(\mathbf{n}, \mathbf{s})) \bullet$ $\quad (\mathbf{shift}(\mathbf{n}, \mathbf{s}))(\mathbf{i}) = \mathbf{s}(\mathbf{i} - \mathbf{n}))$

Recall that the sequence $\langle \mathbf{a}, \mathbf{b} \rangle$ denotes the mapping $\{ 1 \mapsto \mathbf{a}, 2 \mapsto \mathbf{b} \}$. The expression $\mathbf{shift}(3, \langle \mathbf{a}, \mathbf{b} \rangle)$ denotes the mapping $\{(3 + 1) \mapsto \mathbf{a}, (3 + 2) \mapsto \mathbf{b}\}$, which is $\{4 \mapsto \mathbf{a}, 5 \mapsto \mathbf{b}\}$.

Remarks on Z style

We like schema composition: the pipe $()$ and semi-colon $()$ operators. We think of compositions as data-flow diagrams. Pipes connect input and output variables; semi-colon connects initial and final states.

What does this view mean about our style of Z usage? We tend to define intermediate schemas for one of four uses:

- system states,
- state transitions, for use with semi-colon composition,
- data translations, for use with pipes, or
- predicates, for use with logical connectives.

Each of these uses has a different characteristic style. State transitions typically include \exists **State** or Δ **State** for some schema **State**. Data translations typically include both input and output variables, with some axioms defining a relation between them. A predicate typically has input variables, but no output variables. It includes an axiom part (i.e., “the part below the line”) that describes some required properties of the input variables and/or the system state.

Of course, we also define schemas that are merely intended to encapsulate state or define some functions or sets that will be included in later schemas.

Things Left Out of the Model

Our intent is to model significant aspects of the user-program interface to the Synergy file system. Since the Synergy file system interface includes the Unix file system interface as a subset, our model includes aspects of UFS. Our model is not intended to be complete. It neither models all of the user-level operations of the Synergy (or the Unix) file system, nor does it model all the details of every operation included in the model.

Here is a list of some prominent aspects of UFS that we have not included in the SFS model. Some have been left out merely because of time constraints. Some were aspects judged not necessary to produce a useful model. Some were considered unimportant “implementation details” of UFS that show through at the interface level. We have not modeled these aspects, but our model should admit consistent extensions that do include some or all of these aspects.

1. error reporting
2. any resource limits (e.g., size of pathnames, *etc.*)
3. chdir, chroot (at least not yet)
4. chown, chgroup (at least not yet)
5. chmod (at least not yet)
6. mkdir, rmdir
7. mount points to multiple file systems
8. read-only file systems
9. real_uid versus effective_uid
10. real_gid versus effective_gid
11. execution of files (and, hence, the set_uid or set_gid permissions)
12. the directory entries “.” and “..”
13. super-user privileges or restrictions.

14. special files (not even their presence in the file system)
15. soft links (also called symbolic links)
16. file or inode creation/access/modification times
17. blocking or non-blocking accesses
18. some open modes (e.g., `excl`)
19. locking
20. the 4.2 BSD directory manipulation operations: `opendir`, `readdir`, `rewinddir`, `closedir`, `telldir`, `seekdir`.
21. the `stat` system call
22. possible interaction between concurrent reading and writing of a file. (Currently a file is truncated to zero length when it is opened for writing. See the definition of **OpenFT**, page 15.)

3 The Basic File Space

A file is a sequence of bytes.

[**BYTE**]

FILE == seq **BYTE**

The function **padfile** constructs a file containing a sequence of pad bytes to a given length. Usually in Unix, the pad byte is a representation of the number 0.

<p>pad : BYTE padfile : \rightarrow FILE</p>
<p>$\forall n : \bullet \text{padfile}(n) = (\lambda k : 1 \dots n \bullet \text{pad})$</p>

A file identifier is the internal name of a file. The Unix implementation of a file identifier is an inode number.

[**FID**]

The following schema introduces the association of a file identifier with the contents of a file. We specify four transitions on a basic file space: create, destroy, read and write.

<p>BasicFileSpace _____ fcontents : FID FILE</p>
--

NewFID produces a file identifier not currently in use (i.e., a “new file identifier”) as its output. Its output will be piped into other schemas that take a file identifier as input.

NewFid
BasicFileSpace
fid! : FID
fid! \notin dom fcontents

FileLength takes a file identifier as input, and produces an output whose value is the length of the file's contents. The output is named **offset**.

FileLength
BasicFileSpace
fid? : FID
offset! :
fid? \in dom fcontents
offset! = #(fcontents(fid?))

Create

File creation is a state transition on a **BasicFileSpace**. It takes a currently-unused file identifier as input, and causes it to be bound to an empty file in the file space.

CreateBFS
ΔBasicFileSpace
fid? : FID
fid? \notin dom fcontents
fcontents' = fcontents \oplus {fid? \mapsto $\langle \rangle$}

Destroy

When a file is destroyed, the association of a file identifier with the file's contents is removed from the file space.¹

DestroyBFS
ΔBasicFileSpace
fid? : FID
fid? \in dom fcontents
fcontents' = {fid?} fcontents

¹In a typical implementation the disk blocks belonging to the file remain on disk. So removing the association of **fid** \mapsto **file** reflects that the file is no longer accessible via the file system. Our model does not include the underlying storage medium (e.g., a disk system).

Read

The input values of a read operation are a file identifier, an offset where reading begins, and a length of data to be read. The output value is the sequence of bytes that occurs at the given offset in the file. The length of this sequence is the requested length or the remaining length to the end of the file, whichever is smaller. Thus, the length of the resulting sequence depends on the amount of data in the file.

ReadBFS
Ξ BasicFileSpace fid? : FID offset? , length? : data! : FILE
fid? \in dom fcontents data! = (1 .. length?) nthtail (offset? , fcontents (fid?))

Write

The inputs to a write operation are a file identifier, an offset at which writing begins, and a sequence of bytes to be written. The new sequence of bytes replaces existing bytes in the designated range. The file length is increased if necessary, and if the offset is longer than the initial file length pad characters are inserted. ²

WriteBFS
Δ BasicFileSpace fid? : FID offset? : data? : FILE
fid? \in dom fcontents fcontents' = fcontents \oplus { fid? \mapsto padfile (offset?) \oplus fcontents (fid?) \oplus shift (offset? , data?)}

4 The File Space

We introduce the notions of pathname and directory. A pathname is a sequence of syllables. A directory is a mapping from syllables to file identifiers.

[SYLLABLE]

PATHNAME == seq SYLLABLE

DIRECTORY == SYLLABLE FID

²These pad characters correspond to the “holes” that Unix permits within files. These holes are read as pad characters (normally 0). Our model does not explicitly include holes, but does permit writing passed the existing file length, padding the file to the required offset before writing. Thus **read** will see the pad characters in the “hole”, as in the Unix implementation.

We hypothesize the existence of a partial function **ParseDir** which creates an object of type **DIRECTORY** from a file (i.e., makes a file understandable as a directory). The total function **UnParseDir** performs a mapping from directories back to files (in order to store directories as files). This mechanism allows us to represent directories as files in system state, but use a convenient representation when performing “directory operations.” **UnParseDir** is a total function because every directory object (i.e., mapping from name to file identifier) can be represented as a file (i.e., a sequence of bytes).

ParseDir : FILE DIRECTORY UnParseDir : DIRECTORY → FILE
$\text{ran UnParseDir} \subseteq \text{dom ParseDir}$
$\forall \mathbf{d} : \mathbf{DIRECTORY}; \mathbf{s} : \mathbf{SYLLABLE} \bullet$ $\text{dom}(\text{ParseDir}(\text{UnParseDir}(\mathbf{d}))) = \text{dom } \mathbf{d} \wedge$ $(\text{ParseDir}(\text{UnParseDir}(\mathbf{d}))) (\mathbf{s}) = \mathbf{d}(\mathbf{s})$
$\text{dom}(\text{ParseDir}(\langle \rangle)) = \emptyset$

A file may be one of two types, **regular** or **directory**.³ We apply the function **ParseDir** only to files of type **directory**.

FTYPE ::= **regular** | **directory**

As a convenient abstraction, we define the partial function **DirContents** to map a **fid** identifying a directory file to the corresponding **DIRECTORY** value.

BasicFileSpace DirContents : FID DIRECTORY
$\forall \mathbf{fid} : \mathbf{FID} \mid \mathbf{fid} \in \text{dom } \mathbf{fcontents}$ $\wedge \mathbf{fcontents}(\mathbf{fid}) \in \text{dom } \text{ParseDir} \bullet$ $\text{DirContents}(\mathbf{fid}) = \text{ParseDir}(\mathbf{fcontents}(\mathbf{fid}))$

The schema **FileSpace** introduces directories into the basic file space. Each file identifier has a file type. Files of type **directory** are interpreted as directories with the function **ParseDir**. The range of each directory file contains file identifiers occurring in the file space. There is a root file identifier that points to the top of the file tree.

The requirements on a file space are that:

- each file has a file type (i.e., **regular** or **directory**),
- the root file identifier names a directory file in the file space,
- each directory file can be interpreted as a directory structure,
- each directory provides names only for existing files,

³We do not include the Unix concept of special files in our model.

- a directory that is pointed to by no other directory must either be empty or be the root directory.⁴

This last requirement restricts the graph of directories to be a tree structures. We have not included specifications for the self and parent directory entries within a directory (i.e., “.” and “..” in the Unix file system).⁵

FileSpace
BasicFileSpace fctype : FID FTYPE Rootfid : FID
$\text{dom fctype} = \text{dom fcontents}$ Rootfid $\in \text{dom fcontents} \wedge \text{fctype}(\text{Rootfid}) = \text{directory}$ $\forall \text{fid} : \text{dom fcontents} \bullet$ $\text{fctype}(\text{fid}) = \text{directory} \Rightarrow$ $\text{fcontents}(\text{fid}) \in \text{dom ParseDir}$ $\forall \text{fid} : \text{dom fcontents} \bullet$ $\text{fctype}(\text{fid}) = \text{directory} \Rightarrow$ $\text{ran}(\text{DirContents}(\text{fid})) \subseteq \text{dom fcontents}$ $\wedge (\text{parents} == \{\text{parent_fid} : \text{dom fcontents}; \text{name} : \text{SYLLABLE} \mid$ $\text{fctype}(\text{parent_fid}) = \text{directory}$ $\wedge \text{DirContents}(\text{parent_fid})(\text{name}) = \text{fid} \bullet$ $\text{parent_fid}\} \bullet$ $(\#\text{parents} = 0)$ $\Leftrightarrow (\text{fid} = \text{Rootfid} \vee \text{DirContents}(\text{fid}) = \emptyset)$

*Note the property $\text{dom fctype} = \text{dom fcontents}$. This is asserted wherever the **FileSpace** schema is inherited. Note further that $\Delta\text{FileSpace}$ asserts this property for both the initial and the final state variables. This means that we must be careful that any Create and Destroy operations that include $\Delta\text{FileSpace}$ must update both **fctype** and **fcontents**.*

*Remark
on the
Z*

We state that empty directories and regular files can be disconnected from the file space. Both regular files and directories may have multiple links (i.e., $\#\text{parents} > 1$ in the definition of **FileSpace**).⁶

A file is named by a pathname and a starting directory. The pathname identifies the file relative to the starting directory.

The set **legal_pathnames** in a **FileSpace** contains the pair **(startfid, pn)** if and only if **pn** is a legal pathname relative to **startfid** in the file space.

⁴Unix file systems generally depend on a stronger version of this requirement: that a directory file is pointed to by at most one directory entry. Some Unix file system implementations permit this invariant to be violated (e.g., by a privileged user). However, they may not function properly when it is violated.

⁵The special directory entry “.” semi-violates the requirement that the directory structure be tree-structured.

⁶Most Unix systems only permit privileged processes to make new links to existing directories. It is not clear that it is every desirable to have multiple links to directories. We have tried merely to model what exists.

A pathname ($\text{syll}_1, \text{syll}_2, \dots, \text{syll}_k$) starting from a particular directory is legal in a file space either if it is the empty sequence or if the front of the pathname (i.e., excluding the last syllable) is a legal pathname that names a directory and that directory maps the last syllable to a file identifier. When pathname is legal, the function pathname_fid is the file identifier that pathname identifies.

The set directory_fids is the set of fids in the file space that identify directory files. (Specifically, it is the set of fid that the function ftype maps to directory .) This will be convenient, so that we can quantify over the directories in a file space.

LegalPathnames FileSpace $\text{legal_pathnames} : (\text{FID} \times \text{PATHNAME})$ $\text{pathname_fid} : \text{FID} \times \text{PATHNAME} \rightarrow \text{FID}$ $\text{directory_fids} : \text{FID}$
$\text{directory_fids} = \text{ftype}^{-1}\{\text{directory}\}$ $\text{dom pathname_fid} \subseteq \text{directory_fids} \times \text{PATHNAME}$ $\forall \text{fid} : \text{FID} \mid \text{ftype}(\text{fid}) = \text{directory} \bullet$ $\text{pathname_fid}(\text{fid}, \langle \rangle) = \text{fid}$ $\forall \text{fid} : \text{FID}; \text{pn} : \text{PATHNAME}; \text{name} : \text{SYLLABLE} \mid$ $((\text{fid}, \text{pn}) \in \text{dom pathname_fid}$ $\wedge \text{ftype}(\text{pathname_fid}(\text{fid}, \text{pn})) = \text{directory}$ $\wedge \text{name} \in \text{dom}(\text{DirContents}(\text{pathname_fid}(\text{fid}, \text{pn})))) \bullet$ $\text{pathname_fid}(\text{fid}, \text{pn} \langle \text{name} \rangle) = \text{DirContents}(\text{pathname_fid}(\text{fid}, \text{pn}))(\text{name})$ $\text{legal_pathnames} = \text{dom pathname_fid}$

The expression $\text{ftype}^{-1}\{\text{directory}\}$ denotes the inverse image of the singleton set containing the *FTYPE* value directory . That is, it represents the set $\{\text{fid} : \text{FID} \mid \text{ftype}(\text{fid}) \in \{\text{directory}\}\}$

Remark
on the
Z

The schema **FileSpace** includes the property that

$$\forall \text{fid} : \text{dom fcontents} \bullet$$

$$\text{ftype}(\text{fid}) = \text{directory} \Rightarrow$$

$$\text{ran}(\text{DirContents}(\text{fid})) \subseteq \text{dom fcontents}$$

From this we can deduce that

$$\text{ran pathname_fid} \subseteq \text{dom fcontents},$$

because the only fids defined to be in the range of pathname_fid are those that are in the range of a parsed directory.

LegalPathname LegalPathnames $\text{startfid?} : \text{FID}$ $\text{pn?} : \text{PATHNAME}$ $(\text{startfid?}, \text{pn?}) \in \text{legal_pathnames}$
--

The schema **LegalPathname** is designed as a predicate schema.

*Remark
on the
Z*

PathnameLookup LegalPathname fid! : FID
fid! = pathname_fid(startfid?, pn?)

The schema **PathnameLookup** is designed as a data translation schema, for use with pipes (). It inherits inputs **startfid?** and **pn?** from **LegalPathname**.

*Remark
on the
Z*

We define **ExistentFile** and **NonExistentFile** as predicate schemas to distinguish whether a (**startfid?**, **pn?**) pair identifies a file in the **FileSpace** state.

ExistentFile FileSpace LegalPathnames startfid? : FID pn? : PATHNAME
(startfid?, pn?) ∈ dom pathname_fid

NonExistentFile FileSpace LegalPathnames startfid? : FID pn? : PATHNAME
(startfid?, pn?) ∉ dom pathname_fid

We define the schemas **RegularFile** and **DirectoryFile** to be predicates indicating that the file identified by a pathname is a **regular** file or a **directory** file, respectively. These predicates will be used later when we must restrict operations that apply only to one kind of file.

RegularFile FileSpace LegalPathnames startfid? : FID pn? : PATHNAME
ftype(pathname_fid(startfid?, pn?)) = regular

DirectoryFile FileSpace LegalPathnames startfid? : FID pn? : PATHNAME ftype(pathname_fid(startfid?, pn?)) = directory
--

EmptyDirectory DirectoryFile dom(DirContents(pathname_fid(startfid?, pn?))) = \emptyset

We specify the successful actions of linking, unlinking and creating a file. Pathnames instead of file identifiers are used to name files. We do not specify read, write or destroy at this level, since pathnames are never used in the interfaces to these operations.

The schema **Rootfid_is_preserved** is used in the schemas that are Δ **FileSpace** but do not change **Rootfid**. At the moment, we do not include any operations that change **Rootfid**. However, the Unix system defines the **chroot** operation, which does change the (apparent) root of the file space. So, we provide for that possibility.

Rootfid_is_preserved Δ FileSpace Rootfid' = Rootfid

Link

The purpose of a link operation is to associate a new pathname with an existing file.

LinkFS allows us to create a legal pathname to a file regardless of whether the file is already present in any directory (e.g., a file created by **CreateBFS**). The **pn** is the name of the new link to be installed, and **fid** is its target. **front(pn)** must name an existing directory, and **last(pn)** must be a name not currently in that directory. **fid** must identify an existing file.

LinkFS
Δ FileSpace LegalPathnames NonExistentFile Rootfid_is_preserved startfid? : FID pn? : PATHNAME fid? : FID
fid? \in dom fcontents #pn? > 0 (startfid?, front(pn?)) \in legal_pathnames (parentfid == pathname_fid(startfid?, front(pn?)) • ftype(parentfid) = directory \wedge last(pn?) \notin dom(DirContents(parentfid)) \wedge (dir == DirContents(parentfid) • (newdir == dir \cup {last(pn?) \mapsto fid?}) • fcontents' = fcontents \oplus {parentfid \mapsto UnParseDir(newdir)}))
ftype' = ftype

We can describe the new directory contents using \cup , rather than \oplus , because we know that last(pn?) is not in the domain of the directory.

Remark
on the
Z

Unlink

The purpose of an unlink operation is to remove a pathname from the file space. Removing a pathname does not remove its associated file. Thus we do not remove the corresponding element from either **fcontents** or **ftype**. (At this level files may be left with no links to them. After we introduce the “file table” and “process table”, we will add “garbage collection” to the file system.)

The input is the (startfid, pn) pair which names an existing file. In a legal pathname, all initial subsequences are also legal pathnames, and name directories. So we know, in particular, that (startfid?, front(pn)) identifies a directory in the schema below.

UnlinkFS Δ FileSpace LegalPathname ExistentFile Rootfid_is_preserved startfid? : FID pn? : PATHNAME <hr/> #pn? > 0 (parentfid == pathname_fid(startfid?, front(pn?)) • (dir == DirContents(parentfid) • (newdir == {last(pn?)} dir • fcontents' = fcontents \oplus {parentfid \mapsto UnParseDir(newdir)})) ftype' = ftype
--

AddFType

The schema **AddFType** is a companion operation to **CreateBFS**. The two operations combine to update **fcontents** and **ftype** when creating a new file.

AddFType ftype, ftype' : FID FTYPE fid? : FID ftype? : FTYPE <hr/> fid? \notin dom ftype ftype' = ftype \oplus {fid? \mapsto ftype?}
--

Note that **FileSpace** includes the property that $\text{dom ftype} = \text{dom fcontents}$, which is not a necessary property of **AddFType**. In fact, **AddFType** is applied precisely when **ftype** does not yet reflect the file type of the new FID **fid?**. So, we must combine **AddFType** with **CreateBFS** so that this invariant holds in both the initial and final states, and then explicitly include Δ **FileSpace** to state that the invariant holds.

Note also that **AddFType** does not state any relation about the initial and final **Rootfid** or **fcontents**. **AddFType** will be used in schemas that include statements of those relations. Thus, **AddFType** is designed as a partial description of a **FileSpace** transition, and is intended to be combined with other Δ **FileSpace** schemas to more completely describe operations on a **FileSpace**.

*Remark
on the
Z*

DeleteFType

The operation **DeleteFType** removes the association of a file identifier with an **FTYPE** from the function **ftype**.

DeleteFType
ftype, ftype' : FID FTYPE
fid? : FID
fid? ∈ dom ftype
ftype' = {fid?} ftype

CreateFS

In a file space, creating a file is a combination of a basic file creation step and a link operation. The **pn** is the name of the new file to be created. The pathname **front(pn)** must name an existing directory, and the syllable **last(pn)** must be a name not currently in that directory. The input **ftype?** gives the type of the newly created file.

In the schema **CreateFS** these inputs are inherited from the schemas **LinkFS** and **AddFType**. (The input **fid?** to **LinkFS** is captured by **NewFid**, and so is not an input of **CreateFS**.) The file space is changed by the addition of an empty file to the file space as specified by **CreateBFS**, the updating of **ftype** to reflect the new file's type, and the creation of a link to the new file as specified by **LinkFS**.

$$\text{CreateFS} \hat{=} \text{NewFid}((\text{CreateBFS} \wedge \text{AddFType} \wedge \text{Rootfid_is_preserved}) \text{LinkFS})$$

Note that CreateBFS and AddFType can be combined using \wedge because they do not overlap in either their outputs or the state that they “change”.

*LinkFS must be added via $,$ because it requires that $\text{fid?} \in \text{dom fcontents}$, and that is only true for the **fcontents'** following **CreateBFS**.*

*Remark
on the
Z*

5 The File Table

A file must be opened before access can occur. Open files are named by elements of the type **OID**.

[OID]

A file may be opened in one of several “modes”, given permission to do so.

ACCESS_MODE ::= ronly | wronly | rdwr | append

The following schema is so named because it models information implemented in the Unix global file table. **ftposn** is the current position for the open file. **ftmode** associates an access mode with the open file. **ftfid** maps an oid to a file identifier.

FileTable
ftposn : OID
ftmode : OID ACCESS_MODE
ftfid : OID FID
$\text{dom } \mathbf{ftposn} = \text{dom } \mathbf{ftmode} = \text{dom } \mathbf{ftfid}$

We postpone saying $\text{ran}(\mathbf{ftfid}) \subseteq \text{dom}(\mathbf{fcontents})$ until later, because the **FileSpace** schema is not part of the **FileTable** schema. This will be introduced in the **BasicAccessSystem** (section 8, page 22).

*Remark
on the
Z*

Define some data translation schemas for **OID**'s.

OidToFid
FileTable
oid? : OID
fid! : FID
$\mathbf{oid?} \in \text{dom } \mathbf{ftfid}$
$\mathbf{fid!} = \mathbf{ftfid}(\mathbf{oid?})$

OidToPosn
FileTable
oid? : OID
offset! :
$\mathbf{oid?} \in \text{dom } \mathbf{ftposn}$
$\mathbf{offset!} = \mathbf{ftposn}(\mathbf{oid?})$

Open

The schema **OpenFT** specifies a successful file open operation. The file to be opened is specified by **fid?**. **Oid!** is a new open file identifier. The input **access_mode?** specifies the access mode to the open file. The initial position of the opened file is zero.

Directories are updated via the **LinkFS** and **UnlinkFS** operations. The restriction that directory files only be opened in **rdonly** mode is not imposed until later (see section 8, The Basic Access System), because **FileTable** does not include the schema **FileSpace**, and so does not include **ftype**.

When a file is opened for writing (but not for appending), the file is truncated to zero length. This corresponds to requiring the **O_TRUNC** option of the Unix **open** call. It would be easy to extend the specification

to include opening for write without the `O_TRUNC` option, as well as some additional `open` options (e.g., create, exclusive create).

However, some options, such as `O_SYNC`, are not so simple to include. The `O_SYNC` option requires that there be no output buffering on write operations (i.e., that the `write` call not complete until the required physical I/O has completed.) However, the file system specification does not describe buffered I/O; it describes the required changes in the state of the file space. Buffered I/O is an optimization strategy in the implementation. Optimizations should not alter the semantics of an operation, only its performance. A commonly accepted requirement for an optimization is that it not change the semantics of a successful operation, but that it may change the semantics of operations that fail (e.g., a different error may be reported than in the unoptimized implementation). However, we observe that the Unix file system admits buffered I/O in its fundamental file I/O model, and so includes options for dealing with the effects of buffering on file operations. We have not included buffered I/O in our specification of SFS.

If a file is open for reading at the time that it is opened for writing, then this truncation can be observed. There are other choices for the meaning of opening a file for writing when it is already open for reading. Other choices may make the file system specification more complex. We have chosen one that seems particularly simple, although it does not model all of the rich, complex behavior of the Unix file system.

OpenFT illustrates a pattern for defining a schema that operates on a portion of its state. **OpenFT** operates on the **FileTable**, but also needs to refer to the state of the **BasicFileSpace**. We include the declaration $\Delta\mathbf{FileTable}$ to indicate the former. But we include **BasicFileSpace** for the latter, rather than $\exists\mathbf{BasicFileSpace}$. We do this because we may combine **OpenFT** with other schemas that perform corresponding operations on the **BasicFileSpace** using the pipe operator $()$. If we included $\exists\mathbf{BasicFileSpace}$ here, then **OpenFT** would include the property that $\mathbf{fcontents}' = \mathbf{fcontents}$. Including that property would prevent us from piping between **OpenFT** and a schema that operates on **BasicFileSpace**.

*Remark
on the
Z*

OpenInternalFT

$\Delta\mathbf{FileTable}$
BasicFileSpace
fid? : **FID**
access_mode? : **ACCESS_MODE**
oid! : **OID**

fid? \in dom **fcontents**
oid! \notin dom **ftfid**
ftfid' = **ftfid** \oplus {**oid!** \mapsto **fid?**}
ftmode' = **ftmode** \oplus {**oid!** \mapsto **access_mode?**}
ftposn' = **ftposn** \oplus {**oid!** \mapsto 0}

TruncateFile

$\Delta\mathbf{BasicFileSpace}$
fid? : **FID**

fcontents' = **fcontents** \oplus {**fid?** \mapsto $\langle \rangle$ }

TruncatingOpen
access_mode? : ACCESS_MODE
access_mode? \in { wronly , rdwr }

MaybeTruncateFile $\hat{=}$ **TruncatingOpen** \Rightarrow **TruncateFile**

OpenFT $\hat{=}$ **OpenInternalFT** \wedge **MaybeTruncateFile**

CloseFT

We specify the operation of destroying a file table entry. This is the opposite of **OpenFT**. Note, that this effect does not always occur when closing a file at the Unix File System level; a file table entry is expunged only when no process references it. Since a file table entry may be shared by two process file tables, or (in Unix, but not in this draft specification) by multiple entries in a single process table.

CloseFT
Δ FileTable
oid? : OID
oid? \in dom ftfid
ftfid' = { oid? } ftfid
ftmode' = { oid? } ftmode
ftposn' = { oid? } ftposn

The property **oid** \in dom **ftfid** is not necessary in this schema. However, it does indicate our intention about the use of this operation. The **CloseFT** operation could be extended to operate on oid's that do not have entries the file table, in which case the file table should remain unchanged⁷.

Seek

The schema **SeekFT** specifies the operation of logically updating an open file's position. No data is read or written. **SeekFT** merely updates the **ftposn** entry in the file table. Thus, you can seek to a position outside of the current file without error.

⁷ The expression {**oid?**} **ftfid** yields **ftfid** if **oid?** \notin dom **ftfid**. So we need only remove the restriction **oid?** \in dom **ftfid** from the current schema to effect this change.

SeekFT Δ FileTable oid? : OID offset? :
oid? \in dom ftfid ftposn' = ftposn \oplus { oid? \mapsto offset? } ftmode' = ftmode ftfid' = ftfid

The schema **IncrPosnFT** increments the open file position by a specified amount (**delta?**). This will be used later to describe updating the file position after a read or write operation.

IncrPosnFT Δ FileTable oid? : OID delta? :
oid? \in dom ftfid ftposn' = ftposn \oplus { oid? \mapsto ftposn(oid?) + delta? } ftmode' = ftmode ftfid' = ftfid

SeekEOF \cong **OidToFidFileLengthSeekFT**

We now define three more schemas to be used as predicates. They allow us to distinguish classes of access modes recorded in the **FileTable**.

Opened For Read

OpenedForRead FileTable oid? : OID
ftmode(oid?) \in { rdonly , rdwr }

Opened For Write

OpenedForWrite FileTable oid? : OID
ftmode(oid?) \in { wronly , rdwr }

Opened For Append

OpenedForAppend FileTable oid? : OID
ftmode(oid?) = append

6 Processes

The system is populated with processes, each identified by an element of the set **PID**. Each process has an associated file descriptor table, a mapping from file descriptors (type **FD**) to open file identifiers. Each process has its own space of file descriptors, through which it can reference file table entries.

[**PID**, **FD**]

ProcessTable pfhtable : PID (FD OID)

FdToOid

The schema **FdToOid** translates a pair of inputs (**pid?** and **fd?**), which identify an entry in a **pfhtable**, into the **OID** that entry points to.

FdToOid ProcessTable pid? : PID fd? : FD oid! : OID
pid? ∈ dom pfhtable fd? ∈ dom(pfhtable(pid?)) oid! = (pfhtable(pid?))(fd?)

FdToFid \cong **FdToOidOidToFid**

There are four operations on the **ProcessTable**:

- add a new file-descriptor entry,
- delete an existing file-descriptor entry,
- fork a new process, which inherits its initial **pfhtable** from its parent, and
- duplicate a **pfhtable** entry within that **pfhtable**.

AddProcessFD

AddProcessFD
Δ ProcessTable pid? : PID oid? : OID fd! : FD
pid? \in dom pfhtable fd! \notin dom(pfhtable (pid?)) (newfhtable == (pfhtable (pid?) \oplus { fd! \mapsto oid? } \bullet pfhtable' = pfhtable \oplus { pid? \mapsto newfhtable })

DeleteProcessFD

DeleteProcessFD
Δ ProcessTable pid? : PID fd? : FD
pid? \in dom pfhtable fd? \in dom(pfhtable (pid?)) (newfhtable == { fd? } (pfhtable (pid?)) \bullet pfhtable' = pfhtable \oplus { pid? \mapsto newfhtable })

ForkProcess

We partially specify the Unix fork operation. A fork creates a new process, which inherits the file descriptor table of its parent process. The input **pid?** identifies the parent process. The output **pid!** identifies the child process.

ForkPT
ΔProcessTable pid? : PID pid! : PID
pid? \in dom pfdtable pid! \notin dom pfdtable pfdtable' = pfdtable \oplus {pid! \mapsto pfdtable(pid?)}

DuplicateFD

DuplicateFD
ΔProcessTable pid? : PID fd? : FD fd! : FD
pid? \in dom pfdtable fd? \in dom(pfdtable(pid?)) fd! \notin dom(pfdtable(pid?)) (oid == (pfdtable(pid?))(fd?) \bullet newfdtable == (pfdtable(pid?)) \oplus {fd! \mapsto oid} \bullet pfdtable' = pfdtable \oplus {pid? \mapsto newfdtable})

7 Garbage Collection

In a system state, a file may be referenced in several ways: it may have one or more links from directories, and it may be open to one or more processes. A file's contents are removed from the system only when all references to it disappear.⁸ Here we specify garbage collection by directly describing the presence or absence of **ProcessTable** and directory references.

ExistsLinkReferenceToFid
FileSpace LegalPathnames fid? : FID
\exists pn : PATHNAME \bullet ((Rootfid, pn) \in legal_pathnames \wedge pathname_fid(Rootfid, pn) = fid?)

⁸In the standard Unix implementation, each inode contains a count of all the links referencing it, and each file table entry contains a count of the file descriptors in the process file-tables that reference it. When the reference count of an inode is zero, there are no remaining link references in the directories of the file system. When the reference count of a file-table entry is zero, there are no remaining file descriptor entries in any process' file-table referring to that file-table entry. The Unix file system "deletes" a file (as opposed to merely removing a directory link to a file) when both of these counts are zero.

ExistsFileTableReferenceToFid
FileTable fid? : FID
$\exists \text{oid} : \text{OID} \mid \text{oid} \in \text{dom ftfid} \bullet \text{ftfid}(\text{oid}) = \text{fid?}$

$$\text{ExistsReferenceToFid} \hat{=} \text{ExistsLinkReferenceToFid} \vee \text{ExistsFileTableReferenceToFid}$$

A file may be destroyed when there are no references to it. **FileSpaceGC** specifies that a file is garbage collected if there are no references to it, otherwise there is no state change. This is a $\Delta \text{FileSpace}$ transition, but it also depends on **FileTable**; so we declare it as $\exists \text{FileTable}$ to indicate that it does not affect the **FileTable** state (and also to state that the **FileTable** variables are preserved as outputs of **FileSpaceGC**).

$$\begin{aligned} \text{FileSpaceGC} \hat{=} & (\neg \text{ExistsReferenceToFid} \Leftrightarrow (\text{DestroyBFS} \wedge \text{DeleteFType} \\ & \wedge \text{Rootfid_is_preserved} \wedge \Delta \text{FileSpace})) \\ & \wedge (\text{ExistsReferenceToFid} \Leftrightarrow \exists \text{FileSpace}) \\ & \wedge \exists \text{FileTable} \end{aligned}$$

FileTableGC specifies that a file table entry is garbage collected if there are no references to it, otherwise there is no state change. This transition does not depend on the **FileSpace** state, so we need not include $\exists \text{FileSpace}$ here.

ExistsReferenceToFileTableEntry
ProcessTable oid? : OID
$\exists \text{pid} : \text{dom pfdtable} \bullet \exists \text{fd} : \text{dom}(\text{pfdtable}(\text{pid})) \bullet$ $(\text{pfdtable}(\text{pid}))(\text{fd}) = \text{oid?}$

$$\begin{aligned} \text{FileTableGC} \hat{=} & (\neg \text{ExistsReferenceToFileTableEntry} \Leftrightarrow \text{CloseFT}) \\ & \wedge (\text{ExistsReferenceToFileTableEntry} \Leftrightarrow \exists \text{FileTable}) \end{aligned}$$

8 The Basic Access System

The state of the Basic Access System includes the file space, the file table and the process table. A **BasicAccessSystem** has several additional properties compared to a **FileSpace**.

- File table entries must point to files in the **FileSpace**.
- Process file descriptor table entries must point to file table entries. (That is, every process **fd** must map to a file table entry.)

- Every file table entry is pointed to by some process file descriptor table entry. (That is, every file table entry is the image of some process' file descriptor.)
- Every file in the **FileSpace** must either be reachable from the root or be currently open.

The second and third properties together state that all **oid**'s accessible from the **pfddtable** occur in the file table, and all **oid**'s in the file table also occur somewhere in the **pfddtable**. This can be stated succinctly in Z:

$$\text{dom } \mathbf{ftfid} = \bigcup \{ \mathbf{pid} : \text{dom } \mathbf{pfddtable} \bullet \text{ran}(\mathbf{pfddtable}(\mathbf{pid})) \}$$

This says that the domain of **ftfid** (i.e., all open file **oid**'s) is the same as the union of the ranges of all process file tables (i.e., all of the **oid**'s reachable from some **fd** in the process table for some **pid**).

These two properties require garbage collection of inaccessible entries in the **FileTable**. The fourth property additionally requires garbage collection of the **FileSpace**. To state it we introduce the function **rooted_pathname_fid**, which is simply **pathname_fid** with the restriction that the starting **FID** be the **Rootfid** of the **FileSpace**.

RootedPathnames LegalPathnames rooted_pathname_fid : PATHNAME FID <hr/> $\forall \mathbf{pn} : \text{PATHNAME} \bullet$ $\quad \mathbf{rooted_pathname_fid}(\mathbf{pn}) = \mathbf{pathname_fid}(\mathbf{Rootfid}, \mathbf{pn})$
--

BasicAccessSystem FileSpace LegalPathnames RootedPathnames FileTable ProcessTable <hr/> $\text{ran } \mathbf{ftfid} \subseteq \text{dom } \mathbf{fcontents}$ $\text{dom } \mathbf{ftfid} = \bigcup \{ \mathbf{pid} : \text{dom } \mathbf{pfddtable} \bullet \text{ran}(\mathbf{pfddtable}(\mathbf{pid})) \}$ $\forall \mathbf{fid} : \text{dom } \mathbf{fcontents} \bullet$ $\quad \mathbf{fid} \in \text{ran } \mathbf{rooted_pathname_fid} \vee$ $\quad \mathbf{fid} \in \text{ran } \mathbf{ftfid}$
--

Open

Open updates the file table and process table in the basic access system, but not the file space. The file table update is described in the schema **OpenFT**. The process table is updated to map a new file descriptor to a new file table entry. The file must already exist. The **Create** operation is used to create a new file.

Directories are updated only by the **LinkBAS** and **UnlinkBAS** operations. Directory files cannot be changed via the write operation. This restriction is imposed by including the schema **OnlyOpenDirectoriesRDONLY** in **OpenBAS**.

OnlyOpenDirectoriesRDONLY FileSpace fid? : FID access_mode? : ACCESS_MODE
fid? ∈ dom ftype
ftype(fid?) = directory ⇒ access_mode? = ronly

OpenBAS $\hat{=}$ **PathnameLookup**
(OnlyOpenDirectoriesRDONLY \wedge **OpenFT)**
AddProcessFD

*If our model were to stop at the **BAS** level, we would include \exists **FileSpace** in the definition of **OpenBAS**, since an open operation does not change the **FileSpace**. However, we are building toward a model of **UFS** and **SFS**. So we will leave out the \exists **FileSpace**, so that later at the **UFS** level we can pipe from **CreateBAS**, which does change the **FileSpace**, to **OpenBAS**. Similarly, we will not include \exists **FileTable** in the definition of **CreateBAS**.*

*Remark
on the
Z*

Close

CloseBAS $\hat{=}$ (**FdToOid** \wedge **FdToFid**)
(DeleteProcessFD
FileTableGC
(FileSpaceGC \wedge \exists **FileTable))**

*Without the \exists **FileTable**, gives the signature of **CloseBAS** as excluding the **FileTable'** state variables **ftfid'** and **friends**. The maps incoming **ftfid'** to **ftfid**. Even though **FileTableGC** is declared Δ **FileTable**, apparently the ...**FileSpaceGC** eats up the **FileTable'** and, since **FileSpaceGC** isn't Δ **FileTable**, it doesn't produce a new set of "final state" variables.*

*Remark
on the
Z*

Link

The **LinkBAS** schema takes the same logical inputs as **LinkFS**, but the names are different. The existing file (the target of the link) is identified by the input pair (**xstartfid?**, **xpn?**), rather than an **fid**. The input pair (**startfid?**, **newpn?**), inherited from **LinkBAS**, still identifies the new pathname. **LinkBAS**

LinkBAS $\hat{=}$ (**PathnameLookup**[**xstartfid?**/**startfid?**, **xpn?**/**pn?**]
LinkFS)
 \wedge \exists **FileTable** \wedge \exists **ProcessTable**

Unlink

The file to be unlinked is identified by (**startfid?**, **pn?**). That file must exist and be either a regular file or an empty directory.

The prohibition against unlinking non-empty directories simplifies the problem of garbage collecting the file space, because only the file being unlinked is a candidate for garbage collection. If unlinking non-empty directories were permitted, then the entire tree of files under the newly unlinked directory would have to be considered.

$$\text{UnlinkBAS} \hat{=} \text{PathnameLookup} \\ ((\text{RegularFile} \vee \text{EmptyDirectory}) \wedge \text{UnlinkFS FileSpaceGC})$$

Create

The **BasicAccessSystem** added **FileTable** and **ProcessTable** to the **FileSpace**. Since creating a file affects neither the **FileTable** nor the **ProcessTable**, **CreateBAS** need add nothing to **CreateFS**. It takes the same inputs (**startfid?**, **pn?**) and **fctype?** and has signature $\Delta\text{FileSpace}$.

$$\text{CreateBAS} \hat{=} \text{CreateFS}$$

Read

The file to be read is specified by a file descriptor (**FD**). The file must already be opened in one of the access modes that allows a read operation. The open file's position is incremented by the amount of data read.

DataLength
data? : FILE
data! : FILE
delta! :
delta! = # data?
data! = data?

$$\text{ReadBAS} \hat{=} \text{FdToOid} \\ (\text{OidToFid} \wedge \text{OpenedForRead} \wedge \text{OidToPosn}) \\ \text{ReadBFSDataLength}(\text{FdToOidIncrPosnFT})$$

Write

The file to be written is specified by a file descriptor (**FD**). The file must already be opened in one of the access modes that allows a write operation. The offset where writing begins is given by the file table

position, and therefore this input parameter is hidden in the ultimate specification for **WriteBAS**. The file table state is updated to reflect the number of bytes written.

Note that the behavior of writing to a file when it is opened for *write* differs from when it is opened for *append*. When opened for write, the data is written at the current file position (as represented in **ftposn** in **FileTable**). When opened for append, the file position is always set to the end of the file before the data is written. In both cases the file position is incremented by the length of the data written.

The *write* operation is described just below in the **WriteBAS** schema; the *append* operation is described separately in the **AppendBAS** schema.

$$\begin{aligned} \mathbf{WriteBAS} \cong & \mathbf{FdToOid} \\ & (\mathbf{OidToFid} \wedge (\mathbf{OpenedForWrite} \vee \mathbf{OpenedForAppend}) \wedge \mathbf{OidToPosn}) \\ & \mathbf{WriteBFS} \\ & \mathbf{DataLength} \setminus (\mathbf{data!}) \\ & (\mathbf{FdToOidIncrPosnFT}) \end{aligned}$$

*Observe that in **WriteBAS** the schema **DataLength** gets its input **data?** from the outer-most inputs. **WriteBFS** doesn't capture it, because it's not an output being fed down the pipe. But in **ReadBAS**, **DataLength** does capture the **data!** output from **ReadBFS**, and so it must propagate it to the next stage of the pipe. Since we don't want **data!** as an output of **WriteBAS**, we hide it as an output of **DataLength**.*

*Remark
on the
Z*

Append

The append operation is similar to the write operation, but always positions the file at EOF before the data is written.

$$\mathbf{AppendBAS} \cong (\mathbf{FdToOidSeekEOF}) \mathbf{WriteBAS}$$

*Note that the **FdToOid** eats the **fd?** input and produces the **oid?** input required by **SeekEOF**. But the semicolon operator $()$ sees the original **fd?** from the input arguments to the top-level schema-expression, which is required as input to **WriteBAS**.*

*Remark
on the
Z*

Fork

Fork describes the effect on the file system of forking a new process. The new process has its own process file table, which is initially a copy of the file table of its parent process. This operation creates a new *pid* and adds the appropriate entry to the process table. The work is done by **ForkPT**.

$\begin{aligned} & \mathbf{ForkBAS} \\ & \Delta \mathbf{BasicAccessSystem} \\ & \exists \mathbf{FileTable} \\ & \exists \mathbf{FileSpace} \\ & \mathbf{ForkPT} \end{aligned}$
--

DAC_Process_Attributes

```
owner : USER
groups : GROUP
```

The type **DACops** names the operations that to which the discretionary access will apply. Note that we have introduced the **search** operation here. Now that directories have permissions associated with them, a process may be prohibited from searching for a name in a directory. Directory permissions may permit searching for a name in a directory but not reading the directory. That is, a process may be allowed to search for a specific name in a directory, but not read the directory to learn all of the names in it. The **x** file permission is interpreted as the **search** permission for a directory.

DACops ::= open | link | unlink | create | search

DAC

```
fdac : FID DAC_File_Attributes
pdac : PID DAC_Process_Attributes
DAC_current_perms : PID × FID PERMS
DAC_necessary_perms : DACops × ACCESS_MODE PERMS
```

∀ pid : PID; fid : FID •

```
DAC_current_perms(pid, fid) =
  (pdac pid).owner = (fdac fid).owner
  (fdac fid).fperms.owner_perms
  (fdac fid).group ∈ (pdac pid).groups
  (fdac fid).fperms.group_perms
  (fdac fid).fperms.other_perms
```

DAC_necessary_perms(open, ronly) = {r}

DAC_necessary_perms(open, rdwr) = {r, w}

DAC_necessary_perms(open, wronly) = {w}

DAC_necessary_perms(open, append) = {w}

∀ access_mode : ACCESS_MODE •

```
DAC_necessary_perms(create, access_mode) = {x, w} ∧
DAC_necessary_perms(link, access_mode) = {x, w} ∧
DAC_necessary_perms(unlink, access_mode) = {x, w} ∧
DAC_necessary_perms(create, access_mode) = {x, w} ∧
DAC_necessary_perms(search, access_mode) = {x}
```

The schema **DAC** declares the functions **fdac** and **pdac** that map an **fid** or a **pid** to some attributes. The schema **DAC_File_Attributes** (in effect) defines a record structure to represent these attributes (**fowner**, and **fgroup**). **DAC_Process_Attributes** defines an analogous structure for processes.

Remark
on the
Z

The operation **DestroyFileAttrDAC** is never invoked directly. It is only invoked implicitly by **UnlinkOK_UFS**.

DestroyFileAttrDAC
Δ DAC BasicAccessSystem LegalPathname pid? : PID startfid? : FID pn? : PATHNAME
fid == pathname_fid (startfid? , pn?) • fdac' = { fid } fdac
pdac' = pdac

The schema **DACVisiblePathnames** defines what pathnames a process can utilize in DAC file system operations. If a pathname is not “visible” then the user program cannot use that name. Note that visibility is a property of a pathname, not of a file. If a file has several links to it, each link has a different pathname. Some of these pathnames may be visible to a given process and some may not be. Note that if a pathname is in **DACvisible_pathnames**(**pid**, **fid**), then it must also be a legal pathname.

Visible pathnames are described as a function from **PID** \times **FID**. The **FID** represents the starting directory for resolving the pathname. At the user level this will be either the root fid or the fid of the process’ current working directory. If the process does not have search access to a directory, then no pathname starting from that directory will be visible. Thus, a process will not be able to refer to absolute if it does not have search access to the root directory, and will not be able to refer to relative pathnames if it does not have search access to its current working directory.

DACVisiblePathnames
FileSpace LegalPathnames DAC DACvisible_pathnames : PID \times FID \rightarrow (PATHNAME)
\forall pid : PID ; fid : FID ; pn : PATHNAME • $\langle \rangle \in \text{DACvisible_pathnames}(\text{pid}, \text{fid})$ $\Leftrightarrow \text{DAC_necessary_perms}(\text{search}, \text{ronly}) \subseteq \text{DAC_current_perms}(\text{pid}, \text{fid})$ \wedge $(\#\text{pn} > 0$ $\quad \wedge \text{pn} \in \text{DACvisible_pathnames}(\text{pid}, \text{fid}))$ $\Leftrightarrow \text{front}(\text{pn}) \in \text{DACvisible_pathnames}(\text{pid}, \text{fid})$ $\quad \wedge (\text{fid}, \text{pn}) \in \text{legal_pathnames}$ $\quad \wedge \text{DAC_necessary_perms}(\text{search}, \text{ronly})$ $\quad \subseteq \text{DAC_current_perms}(\text{pid}, \text{pathname_fid}(\text{fid}, (\text{front}(\text{pn}))))$

DACVisiblePathname DACVisiblePathnames pid? : PID startfid? : FID pn? : PATHNAME <hr/> pn? ∈ DACvisible_pathnames(pid?, startfid?)

The **DAC** test for each operation is basically

$$\text{DAC_necessary_perms}(\text{Operation}, \text{Access_Mode}) \subseteq \text{DAC_current_perms}(\text{Pid}, \text{Fid})$$

Operations will also require that their pathname inputs be appropriately visible. (E.g., for the link operation, the **front** of the pathname to be created must exist, but the full pathname must not exist.)

*Since we defined **DAC_necessary_perms** to take both a **DACop** and an **access_mode** as arguments, we must associate an **access_mode** with each operation when checking **DAC** permissions. Since **Link** and **Unlink** do not normally have any associated access modes, we adopt the arbitrary convention of associating them with **rdwr** mode.*

*Remark
on the
Z*

10 The Unix File System Interface

The UFS level combines the **BasicAccessSystem** with the **DAC**. In addition, it restricts some operations according to file type. A number of these restrictions appear to be a mechanism to enforce a protocol on file manipulation on processes so that desirable UFS invariants can be enforced. For example, there can only be a single link to a directory file. (This means that the file space underlying the UFS is restricted to being a rooted tree, instead of an arbitrary graph.)

At the UFS level, **link** and **unlink** are only permitted on regular files, not on directories. Only one link to a directory is permitted, the one created when the directory is created. That link to a directory is removed by the **rmdir** operation, which is just like **UnlinkFS**, but requires additionally that the directory be empty. Also, rather than a single **create** operation for both files and directories, **create** is restricted to creating **regular** files and the new operation **mkdir** is used to create **directory** files. The **write** operation is also not permitted on directories; directories are only updated by **link**, **unlink**, **create**, **mkdir**, and **rmdir**.

The “current process” is represented by the current process identifier, **cpid**. Each process has a “working directory”. This will be used as the starting **FID** when resolving “relative pathnames.” The mapping from each **pid** to its working directory is represented in the state variable **pcwd**. Note that there is no constraint that $\text{ran pcwd} \subseteq \text{dom fcontents!}$ There is no mechanism to prevent a process’ from deleting (i.e., unlinking) a directory that is in use as a working directory. Unix permits this, and so we do, as well.¹⁰ Further, user process need not have read/search access to either the root directory or to its current working directory. Of course, in that case the process will not be able to reference absolute or relative pathnames (respectively).

¹⁰If **chmod** were included in our model, we would require that the new working directory exist – and be a directory – in order for the **chmod** operation to be successful.

UFS
BasicAccessSystem
DAC
LegalPathnames
pcwd : PID FID
cpid : PID
dom fdac = dom fcontents
dom pdac = dom pfdtable
dom pcwd = dom pfdtable
ran pcwd \subseteq directory_fids
cpid \in dom pfdtable

The transitions on the controlled access system are those of the basic access system, with additional access control constraints. At this level, pathnames are labeled either relative or absolute.

PATHNAMETYPE ::= relative | absolute

LABELED_PATHNAME == PATHNAMETYPE \times PATHNAME

The functions **lpntype** and **lpnpathname** extract the type and pathname of a labeled pathname, respectively.

lpntype : LABELED_PATHNAME \rightarrow PATHNAMETYPE
lpnpathname : LABELED_PATHNAME \rightarrow PATHNAME
$\forall \text{lpn} : \text{LABELED_PATHNAME}; \text{pnt} : \text{PATHNAMETYPE} \bullet$ $\text{lpntype}(\text{lpn}) = \text{pnt} \Leftrightarrow (\exists \text{pn} : \text{PATHNAME} \bullet (\text{pnt}, \text{pn}) = \text{lpn})$
$\forall \text{lpn} : \text{LABELED_PATHNAME}; \text{pn} : \text{PATHNAME} \bullet$ $\text{lpnpathname}(\text{lpn}) = \text{pn} \Leftrightarrow (\exists \text{pnt} : \text{PATHNAMETYPE} \bullet (\text{pnt}, \text{pn}) = \text{lpn})$

AbsolutePathname

Operations at the UFS level accept input of type **LABELED_PATHNAME**. The **AbsolutePathname** schema translates from a **LABELED_PATHNAME** input to the **(startfid?, pn?)** pair used by the **BasicAccessSystem** operations. **AbsolutePathname** is designed for piping.

AbsolutePathname
UFS
lpn? : LABELED_PATHNAME
pid! : PID
startfid! : FID
pn! : PATHNAME
pid! = cpid
startfid! = lpntype(lpn?) = relative pcwd(cpid) Rootfid
pn! = lpnpathname(lpn?)

Create

The **CreateOK_UFS** operation represents a successful call **Create** to the Unix kernel. We will describe error returns and error codes later.

For a **Create** operation, the process must have **r** and **x** access to the directory in which the file is to be created. **CreateUFS** takes as input the pathname for the new file (**lpn?**). **AbsolutePathname** is used to translate the labeled pathname **lpn?** to **startfid?** and **pn?**, which are the inputs expected by the **BasicAccessSystem** operations. At the UFS level **Create** combines **CreateBAS** and **OpenBAS**, so that it both creates a new file and produces a new **FD** for that file as output.

The permission checking schema, **CreatePermittedDAC**, checks that current process has the proper permissions for the directory in which the new file is to be created.

CreatePermittedDAC
DAC
BasicAccessSystem
DACVisiblePathname
startfid? : FID
pn? : PATHNAME
front(pn?) ∈ DACvisible_pathnames(pid?, startfid?)
(fid == pathname_fid(startfid?, front(pn?)) • DAC_necessary_perms(create, rdwr) ⊆ DAC_current_perms(pid?, fd))

Note that it is possible to create a file with permissions that say you can't write to it. And since create opens the file, too, you get the **fd** back so you **can** write to it while it's open!

The **CreateFileAttrsDAC** operation creates the initial **DAC_File_Attributes** for a newly created file. The file's owner is taken from the power of the process creating the file. The file's group is taken from the group of its parent directory in the BSD style. (This means that we cannot create the root directory this way!) The file's permissions are passed in via the **fperms?** input.

CreateFileAttrsDAC
Δ DAC BasicAccessSystem LegalPathnames pid? : PID startfid? : FID pn? : PATHNAME fperms? : FPERMS
pathname_fid(startfid?, pn?) \notin dom fdac #pn? > 0 (fid == pathname_fid(startfid?, pn?) • (parent_fid == pathname_fid(startfid?, front(pn?)) • (fdac'(fid)).fperms = fperms? \wedge (fdac'(fid)).owner = (pdac(pid?)).owner \wedge (fdac'(fid)).group = (fdac(parent_fid)).group)) fdac = {pathname_fid(startfid?, pn?)} fdac' pdac' = pdac

CreateOK_UFS $\hat{=}$ AbsolutePathname
(CreatePermittedDAC \wedge CreateBAS \wedge OpenBAS
\wedge CreateFileAttrsDAC)
 \wedge Δ **DAC**

Open

The **OpenOK_UFS** operation represents a successful call to **open** in the Unix kernel. We will describe error returns and error codes later.

We have not tried to model all options of the Unix **open** system call. However, we believe that the schemas defined so far are sufficient to do so.

Under Unix, **Open** may either open an existing file (as in lower levels) or create a file. We model this by having **Open** call **Create** if the named file does not exist.

OpenPermittedDAC
DAC BasicAccessSystem DACVisiblePathname startfid? : FID pn? : PATHNAME access_mode? : ACCESS_MODE
(fid == pathname_fid(startfid?, pn?) • DAC_necessary_perms(open, access_mode?) \subseteq DAC_current_perms(pid?, fid))

$$\text{OpenExistingFileOK_UFS} \hat{=} \text{AbsolutePathname}(\text{OpenPermittedDAC} \wedge \text{OpenBAS}) \\ \wedge \exists \text{DAC}$$

Link and Unlink

The operations **LinkOK_UFS** describes the outcome of a successful link operation on a **UFS** state. It takes two arguments of type **LABELED_PATHNAME**, **lpn?** and **xlpn?**. The input **lpn?** is the name of the new line to be created, and must not name an existing file. The input **xlpn?** must name an existing, regular file. Both **xlpn?** and **front(lpn?)** must be pathnames visible to the current process.

LinkPermittedDAC
DAC BasicAccessSystem DACVisiblePathnames pid? : PID startfid?, xstartfid? : FID pn?, xpn? : PATHNAME
front(pn?) \in DACvisible_pathnames(pid?, startfid?) xpn? \in DACvisible_pathnames(pid?, xstartfid?) (fid == pathname_fid(startfid?, front(pn?)) • DAC_necessary_perms(link, rdwr) \subseteq DAC_current_perms(pid?, fid)

$$\text{LinkOK_UFS} \hat{=} (\text{AbsolutePathname} \wedge \\ \text{AbsolutePathname}[\text{xlpn?}/\text{lpn?}, \text{xstartfid!}/\text{startfid!}, \text{xpn!}/\text{pn!}]) \\ (\text{LinkPermittedDAC} \wedge \\ \text{LinkBAS}) \\ \wedge \exists \text{DAC}$$

UnlinkPermittedDAC
DAC BasicAccessSystem LegalPathnames DACVisiblePathname pid? : PID startfid? : FID pn? : PATHNAME
fid == pathname_fid(startfid?, pn?) • DAC_necessary_perms(unlink, rdwr) \subseteq DAC_current_perms(pid?, fid)

$$\begin{aligned} \mathbf{UnlinkOK_UFS} &\hat{=} \mathbf{AbsolutePathname} \\ &(\mathbf{UnlinkPermittedDAC} \wedge \mathbf{RegularFile} \wedge \mathbf{UnlinkBAS}) \\ &\wedge \Delta \mathbf{DAC} \end{aligned}$$

The UFS requirement that

$$\text{dom } \mathbf{fdac} = \text{dom } \mathbf{fcontents}$$

implies that if **UnlinkOK_UFS** results in the file being GC'd, then the DAC file attributes (**fdac(fid)**) must be “GC'd” as well. **DestroyFileAttrDAC** describes removing **fid** from the domain of **fdac**, but **DestroyFileAttrDAC** is never invoked explicitly.

Read, Write, Close

Read and write operations are permitted on open files. The intended DAC constraint is that a process can only read or write files to which it has permission. But the **open** operation/kernel call acts as a “gatekeeper” for read, write, and close. So once the **open** operation has checked permission, read and write need only check that the file was opened in an appropriate access mode. The operation **WriteOK_UFS** combines the two **BasicAccessSystem** operations **WriteBAS** and **AppendBAS**, as the Unix kernel call **write** does.

If the permissions associated with a file are constant, then a process will only be permitted to read or write a file to which it has the appropriate permission. However, if file permissions can change dynamically, then the file system no longer has this property. Rather, a weaker condition holds: a process will only be permitted to read or write a file to which it or an ancestor had appropriate permission when the **Open** operation was performed.

$$\begin{aligned} \mathbf{ReadOK_UFS} &\hat{=} \mathbf{UFS} \wedge \mathbf{ReadBAS}[\text{cpid/pid?}] \\ &\wedge \exists \mathbf{DAC} \end{aligned}$$

$$\begin{aligned} \mathbf{CloseOK_UFS} &\hat{=} \mathbf{UFS} \wedge \mathbf{CloseBAS}[\text{cpid/pid?}] \\ &\wedge \exists \mathbf{DAC} \end{aligned}$$

$$\begin{aligned} \mathbf{WriteOK_UFS} &\hat{=} \mathbf{UFS} \\ &\wedge (((\mathbf{FdToOid}[\text{cpid/pid?}]\mathbf{OpenedForWrite}) \wedge \mathbf{WriteBAS}[\text{cpid/pid?}]) \\ &\quad \vee ((\mathbf{FdToOid}[\text{cpid/pid?}]\mathbf{OpenedForAppend}) \wedge \mathbf{AppendBAS}[\text{cpid/pid?}])) \\ &\wedge \exists \mathbf{DAC} \end{aligned}$$

We use the form **predicate_i ∧ state_transition_i** to indicate that the **ith** *state_transition* only applies when the **ith** *predicate* holds. We use this form, rather than phrasing it as **predicate_i ⇒ state_transition_i**, because the implication is vacuously true if the predicate is not satisfied. Thus the schema is vacuously satisfied if none of the predicates are satisfied. When we add error reporting and return codes to our model, we want to be sure that the schemas modeling successful transitions are not also satisfied (even vacuously) under conditions that should report an error.

*Remark
on the
Z*

11 Mandatory Access Controls

We treat the specification of the MAC layer analogously to the specification of the DAC layer. Both access control layers have an access control policy that controls an access mechanism through an interface-language of “permissions”. The DAC and MAC policies take consider different attributes of processes and files, and they deal with different languages of permissions, but the same abstract structure describes both.

The Synergy MAC layer uses *extended permissions* to provide finer-grained control of file access. Each file and each process has an associated *security context*, which identifies the security-related attributes needed to make a policy decision. The MAC policy produces a set of extended permissions describing permitted access based on the security contexts of the user requesting access and the file being accessed.

The MAC operations and permissions described here are based on informal discussions of the Synergy file system in Fall 1994. These specifications do not reflect the full semantics of Synergy prototype.

```
MAC_PERM ::= read_perm | write_perm | exec_perm | append_perm |
           create_perm | delete_perm | search_perm |
           get_context_perm | set_context_perm
```

```
MAC_PERMS == (MAC_PERM)
```

```
[SECCONTEXT]
```

As in the DAC layer, we define schemas to encapsulate the file and process attributes relevant to the policy decision.

```
MAC_File_Attributes
context : SECCONTEXT
```

```
MAC_Process_Attributes
context : SECCONTEXT
```

As in the DAC layer, we define the set of operations a process may perform on a file. The policy decision is based on the MAC attributes of the process and the file, and the operation being requested.

The SFS open operation accepts several modifiers, called modes. The description of an operation is a pair, a **MACops** and an **EXT_MODE**.

```
MACops ::= mac_open | mac_link | mac_unlink | mac_create |
           mac_search | mac_append | mac_delete |
           mac_get_context | mac_set_context
```

```
EXT_MODE ::= xm_read | xm_write | xm_readwrite | xm_append |
            xm_truncate | xm_excl
```

We define the MAC schema quite analogously to the earlier DAC schema. The two functions **fmac** and **pmac** map files and processes to their MAC-related attributes. The function **MAC_current_perms** maps a **pid** and an **fid** to the set of extended permissions describing how that process is permitted to operate on that file. The definition of that set of permissions is encapsulated in the function **MAC_policy_perms**. This policy computation considers only the relevant MAC attributes of the process and file in question (i.e., their **SECCONTEXTs**). We provide only the signature of **MAC_policy_perms**, leaving the actual policy unspecified. Finally, **MAC_necessary_perms** produces the set of permissions required for an operation to complete without error.

<pre> MAC fmac : FID MAC_File_Attributes pmac : PID MAC_Process_Attributes MAC_current_perms : PID × FID MAC_PERMS MAC_necessary_perms : MACOps × EXT_MODE MAC_PERMS MAC_policy_perms : (SECCONTEXT × SECCONTEXT) MAC_PERMS ∀ pid : PID; fid : FID • MAC_current_perms(pid, fid) = MAC_policy_perms((pmac pid).context, (fmac fid).context) MAC_necessary_perms(mac_open, xm_read) = {read_perm} MAC_necessary_perms(mac_open, xm_readwrite) = {read_perm, write_perm} MAC_necessary_perms(mac_open, xm_write) = {write_perm} MAC_necessary_perms(mac_open, xm_append) = {append_perm} ∀ ext_mode : EXT_MODE • MAC_necessary_perms(mac_search, ext_mode) = {search_perm} ∧ MAC_necessary_perms(mac_create, ext_mode) = {create_perm} ∧ MAC_necessary_perms(mac_delete, ext_mode) = {delete_perm} ∧ MAC_necessary_perms(mac_get_context, ext_mode) = {get_context_perm} ∧ MAC_necessary_perms(mac_set_context, ext_mode) = {set_context_perm} </pre>

We could have used a different decomposition for these security-related system components. For example, we might define one schema for **fmac** and **pmac**, which represent file and process attributes that normally change during system operation, and a separate schema for **MAC_policy_perms** and **MAC_necessary_perms**, which generally change less often, if at all, during system operation. That way we distinguish the policy and enforcement tests from dynamic process and file attributes. So then even operations that may dynamically alter file or process attributes clearly leave the policy unchanged.

Analogously to the DAC layer, we define the notion of a pathname being visible to a process. We do this with two schemas. **MACVisiblePathnames** introduces the function **MAC_visible_pathnames**, which maps a process and a starting **fid** to the set of pathnames that process may “see” in a given file system state. The schema **MACVisiblePathname** takes inputs **pid?**, and (**startfid?**, **pn?**) and asserts the property that **pn?** is a member of the process’s **Mac_visible_pathnames**.

We may want to include both DAC and MAC access controls at the SFS level. This is easily done by including the both relevant MAC schema and the relevant DAC schema in the SFS definition. For example, if we want to assure that a pathname is visible under both the MAC and DAC policies, we might include both **MACVisiblePathname** and **DACVisiblePathname** in a schema definition.

MACVisiblePathnames FileSpace LegalPathnames MAC MACvisible_pathnames : $PID \times FID \rightarrow (PATHNAME)$
$\forall pid : PID; fid : FID; pn : PATHNAME \bullet$ $\langle \rangle \in MACvisible_pathnames(pid, fid)$ \wedge $(\#pn > 0$ $\quad \wedge pn \in MACvisible_pathnames(pid, fid)$ $\Leftrightarrow front(pn) \in MACvisible_pathnames(pid, fid)$ $\quad \wedge (fid, pn) \in legal_pathnames$ $\quad \wedge MAC_necessary_perms(mac_search, xm_read)$ $\quad \subseteq MAC_current_perms(pid, pathname_fid(fid, front(pn))))$

MACVisiblePathname MACVisiblePathnames pid? : PID startfid? : FID pn? : $PATHNAME$
$pn? \in MACvisible_pathnames(pid?, startfid?)$

12 The Synergy File System Interface

This section provides the *form* for describing the SFS layer. The specific details of the operations and MAC tests are merely meant to be illustrative.

SFS UFS MAC
$dom\ fmac = dom\ fcontents$ $dom\ pmac = dom\ pfdtable$

Create

Note that **CreatePermittedMAC** checks the permissions of the directory in which the new file would reside.

CreatePermittedMAC
MAC BasicAccessSystem MACVisiblePathname startfid? : FID pn? : PATHNAME
front(pn?) ∈ MACvisible_pathnames(pid?, startfid?) (fid == pathname_fid(startfid?, front(pn?)) • MAC_necessary_perms(mac_create, xm_readwrite) ⊆ MAC_current_perms(pid?, fid))

The operation **CreateFileAttrsMAC** creates the initial MAC file attributes for a new file. It will probably share with (or perhaps implement) the definition of the **SetContext** operation. However, since it is used here as a part of the **Create** operation, it need not perform any MAC checks.

CreateFileAttrsMAC
Δ MAC LegalPathnames secontext? : SECCONTEXT pid? : PID startfid? : FID pn? : PATHNAME
pathname_fid(startfid?, pn?) ∉ dom fmac #pn? > 0 fmac = {pathname_fid(startfid?, pn?)} fmac' (fmac'(pathname_fid(startfid?, pn?))).context = secontext? pmac' = pmac

CreateFileAttrsMAC

We conjoin the predicate properties of **CreatePermittedMAC** and **CreateOK_UFS**. So the the overall **CreateOK_SFS** operation is only defined when all of their properties are satisfied (e.g., both MAC and DAC tests).

$$\begin{aligned} \text{CreateOK_SFS} \hat{=} & (\text{AbsolutePathname} \\ & (\text{CreatePermittedMAC} \wedge \text{CreateFileAttrsMAC})) \\ & \wedge \text{CreateOK_UFS} \end{aligned}$$

Open

OpenPermittedMAC MAC BasicAccessSystem MACVisiblePathname startfid? : FID pn? : PATHNAME xmode? : EXT_MODE <hr/> (fid == pathname_fid(startfid?, pn?) • MAC_necessary_perms(mac_open, xmode?) ⊆ MAC_current_perms(pid?, fid))

OpenExistingFileOK_SFS $\hat{=}$ (AbsolutePathnameOpenPermittedMAC)
 \wedge OpenExistingFileOK_UFS
 \wedge \exists MAC

Link

LinkPermittedMAC MAC BasicAccessSystem MACVisiblePathnames pid? : PID startfid?, xstartfid? : FID pn?, xpn? : PATHNAME <hr/> front(pn?) \in MACvisible_pathnames(pid?, startfid?) xpn? \in MACvisible_pathnames(pid?, xstartfid?) (fid == pathname_fid(startfid?, front(pn?)) • MAC_necessary_perms(mac_link, xm_readwrite) ⊆ MAC_current_perms(pid?, fid))

LinkOK_SFS $\hat{=}$ (AbsolutePathname \wedge
AbsolutePathname[xlpn?/lpn?, xstartfid!/startfid!, xpn!/pn!])
(LinkPermittedMAC
 \wedge LinkOK_UFS)
 \wedge \exists MAC

Unlink

UnlinkPermittedMAC
MAC BasicAccessSystem LegalPathnames MACVisiblePathname pid? : PID startfid? : FID pn? : PATHNAME
$\text{fid} == \text{pathname_fid}(\text{startfid?}, \text{front}(\text{pn?})) \bullet$ $\text{MAC_necessary_perms}(\text{mac_unlink}, \text{xm_readwrite})$ $\subseteq \text{MAC_current_perms}(\text{pid?}, \text{fid})$

The **Unlink_SFS** operation removes an entry from the directory (if MAC and DAC permissions allow it). The MAC file attributes associated with the file are removed when the file is GC'd from the file space (see **FileSpaceGC**, **CloseBAS**, and **UnlinkBAS**). This follows from the requirement that in all **SFS** states the domain of **fmac** is the same as the domain of **fcontents**. File attributes are associated with files, not with directory entries. Hence they are not necessarily affected each time a directory entry is deleted.

$$\text{UnlinkOK_SFS} \hat{=} (\text{AbsolutePathname} \\ (\text{UnlinkPermittedMAC} \wedge \text{UnlinkOK_UFS})) \\ \wedge \Delta\text{MAC}$$

Read, Write, Close

$$\text{ReadOK_SFS} \hat{=} \text{ReadOK_UFS}$$

$$\text{WriteOK_SFS} \hat{=} \text{WriteOK_UFS}$$

$$\text{CloseOK_SFS} \hat{=} \text{CloseOK_UFS}$$

Get_Context

GetContextPermittedMAC
MAC LegalPathname MACVisiblePathname pid? : PID startfid? : FID pn? : PATHNAME
$(\text{fid} == \text{pathname_fid}(\text{startfid?}, \text{pn?})) \bullet$ $\text{MAC_necessary_perms}(\text{mac_get_context}, \text{xm_readwrite})$ $\subseteq \text{MAC_current_perms}(\text{pid?}, \text{fid})$

GetContext MAC LegalPathname fid? : FID secontext! : SECCONTEXT <hr/> secontext! = (fmac(fid?)).context
--

GetContextOK_SFS $\hat{=}$ **AbsolutePathname**
(GetContextPermittedMAC \wedge
(PathnameLookupGetContext))

13 Reactions to the Z notation

We are both impressed with the simplicity of simple specifications in Z, and with the complexity of anything larger.

1. Z encourages defining a system state in terms a large number of state variables. It provides little mechanism for structuring the state.
2. Standard Z usage requires that a schema representing a state transition specify the new value of all state variables.

That is, the inclusion of ΔS in a schema representing a state transition on the system state described in schema **S** introduces all of the variables in **S'** (i.e., all of the primed state-variables). If the value any of these primed variables are not specified in the state-transition schema, then those values are completely unconstrained and completely unrelated to the corresponding initial state variable (i.e., the corresponding unprimed variable). Since there tend to be many state variables, this mechanism is error prone.

3. Many simple errors in Z specifications prove difficult to uncover.
 - (a) Use of redundant input/output variables in schemas.
If you don't include redundant variable declarations in a schema that inherits them from another schema, then the very variables of interest to your schema are not made explicit in it.
If you do include redundant variable declarations, there is no protection from misspelling variables names.
 - (b) Sometimes primed variables are left unspecified in a schema, because that schema will later be combined with other schemas that do specify that variable. But there is neither a mechanism to declare this intent, nor an easy mechanism to detect which variables are left unspecified. This problem is particularly noticeable as the number state variables gets large, and you attempt to combine schemas (via schema inclusion or composition).

Type checking is a very useful means of finding first-order errors in a Z specification. But it by no means enough.

We have not investigated the “object oriented” variations of Z. Perhaps by applying the structuring principles of object-oriented programming to Z specifications, larger-sized specifications will become more manageable.

As part of this effort an executable model of this specification was written using the ACL2 system. Rephrasing the specification in ACL2 exposed some minor errors. Getting the model accepted by the ACL2 system amounts to a proof that the model is consistent.

14 Errors Uncovered by ACL2 Model

Here is a summary of errors detected during the development of this Z specification by building an ACL2 model of the specification. The ACL2 model only went as far as the DAC layer. (These errors have been corrected, and so do not appear in the Z specification in this report.)

1. *Input name mismatch.* In the schema **LinkFS** the schema **RegularFile** was included “above the line.” **RegularFile** was intended to apply to the existing file, but **RegularFile** relates to inputs **startfid?** and **pn?**, whereas the existing file input for **LinkFS** is **fid?**.
2. *Contradiction.* The schemas **RegularFile** and **NonExistentFile** are contradictory; that is, they denote disjoint sets of bindings, and so no bindings satisfy both schemas. Since **LinkFS** included them both, no bindings would have satisfied **LinkFS**. Of course, **RegularFile** was intended to apply to the existing file, and **NonExistentFile** was intended to apply to the new pathname for that file.
3. *Incorrect specification.* The schema **OpenFT** incorrectly truncated files when opening them for writing, even if the file was being opened in append mode.
4. *Use of undefined expressions.* Zed permits the use of expressions which mention a function applied to arguments not in the domain of that function. Zed says such expressions are meaningless. However, it is not always obvious that the author meant to write a meaningless expression.

In the schema **OnlyOpenDirectoryRDONLY** the expression **ftype(fid?)** is used without requiring **fid? ∈ dom ftype**. Thus, sometimes the schema’s predicate is meaningless. We may know (or believe) that the context of use of always requires **fid? ∈ dom ftype**, because most other operations require it. ACL2 requires that functions only be applied to arguments in their domains. So this casual style was not permitted.

5. *Framing error.* The **LinkBAS** operation did not specify that state variables **ftype** and **RootFid** remained unchanged.
6. *Framing error.* The **CreateBAS** operation did not specify that state variables **FileTable** and **ProcessTable** remained unchanged. However, the real error was that **CreateBAS** should not have been **ΔBasicAccessSystem** at all!
7. *Importing error.* The schema **LinkFS** imported the schema **LegalPathname**, rather than the schema **LegalPathnames**. The former specifies that the pathname identified by (**startfid?**, **pn?**) names an existing file. The schema **LegalPathnames** just introduces the **legal_pathnames** (the set of legal pathnames) as a schema variable. The properties “below the line” in **LinkFS** then related its inputs to the set **legal_pathnames**.

References

- [Bach, 1986] Maurice J. Bach.
The Design of the Unix Operating System.
Prentice-Hall, Englewood Cliffs, New Jersey, 07632, 1986.
- [CSRG, 1986] CSRG.
Unix User's Reference Manual.
Computer Science Division, Dept. of Electrical Engineering and Computer Science, University of California, Berkeley, California, 4.3 berkeley software distribution edition, 1986.
- [Leffler *et al.*, 1989] S.J. Leffler, M.K. McKusick, M.J. Karels, and J.S. Quarterman.
The Design and Implementation of the 4.3BSD Unix Operating System.
Addison-Wesley, 1989.
- [Morgan and Sufrin, 1987] Carol Morgan and Bernard Sufrin.
Specification for the unix filing system.
In Ian Hayes, editor, *Specification Case Studies*, pages 91–140. Prentice-Hall, 1987.
- [Spivey, 1989] J.M. Spivey.
The Z Notation: A Reference Manual.
Prentice Hall, 1989.
- [Stevens, 1992] W. Richard Stevens.
Advanced Programming in the UNIX Environment.
Addison-Wesley, 1992.

Index

- absolute
 - PATHNAMETYPE constant, 31
- AbsolutePathname, 32
- ACCESS_MODE, 14
- AddFType, 13
- AddProcessFD, 20
- append
 - ACCESS_MODE constant, 14
- AppendBAS, 26
- BasicAccessSystem, 23
- BasicFileSpace, 4
- BYTE, 4
- CloseBAS, 24
- CloseFT, 17
- CloseOK_SFS, 41
- CloseOK_UFS, 35
- constants
 - absolute, 31
 - append, 14
 - create, 28
 - link, 28
 - open, 28
 - r (PERM), 27
 - rdonly, 14
 - rdwr, 14
 - regular (FTYPE), 7
 - relative, 31
 - search, 28
 - unlink, 28
 - w (PERM), 27
 - wronly, 14
 - x (PERM), 27
- create
 - DACops constant, 28
- CreateBAS, 25
- CreateBFS, 5
- CreateFileAttrsDAC, 33
- CreateFS, 14
- CreateOK_SFS, 40
- CreateOK_UFS, 33
- CreatePermittedDAC, 32
- CreatePermittedMAC, 39
- DAC, 28
- DAC_File_Attributes, 28
- DAC_Process_Attributes, 28
- DACops, 28
- DACVisiblePathname, 30
- DACVisiblePathnames, 30
- data translations, 3
- DataLength, 25
- DeleteFType, 14
- DeleteProcessFD, 20
- DestroyBFS, 5
- DestroyFileAttrDAC, 29
- DirContents, 7
- DIRECTORY, 7
- directory_fids, 9
- DirectoryFile, 11
- DuplicateFD, 21
- EmptyDirectory, 11
- ExistentFile, 10
- ExistsFileTableReferenceToFid, 22
- ExistsLinkReferenceToFid, 22
- ExistsReferenceToFid, 22
- ExistsReferenceToFileTableEntry, 22
- EXT_MODE, 37
- FD, 19
- FdToFid, 19
- FdToOid, 19
- FID, 4
- FILE, 4
- FileLength, 5
- FileSpace, 8
- FileSpaceGC, 22
- FileTable, 15
- FileTableGC, 22
- ForkBAS, 27
- ForkPT, 21
- FPERMS, 27
- free types
 - BYTE, 4
 - DIRECTORY, 7
 - FD, 19
 - FID, 4
 - FILE, 4
 - FTYPE, 7
 - GROUP, 27
 - OID, 14
 - PATHNAME, 7

- PID, 19
- SYLLABLE, 7
- USER, 27
- FTYPE, 7
- functions
 - DAC_current_perms, 28
 - DAC_necessary_perms, 28
 - DACvisible_Pathnames, 30
 - DirContents, 7
 - directory_fids, 9
 - fcontents, 4
 - fdac, 28
 - fmac, 37
 - ftfid, 15
 - ftmode, 15
 - ftposn, 15
 - legal_pathname, 9
 - lpnpathname, 31
 - lpntype, 31
 - MAC_current_perms, 37
 - MAC_necessary_perms, 37
 - MAC_policy_perms, 37
 - MACVisible_Pathnames, 38
 - padfile, 4
 - ParseDir, 7
 - pathname_fid, 9
 - pdac, 28
 - pfhtable
 - state variable, 19
 - pmac, 37
 - rooted_pathname_fid, 23
 - UnParseDir, 7
- GetContext, 42
- GetContextOK_SFS, 42
- GetContextPermittedMAC, 42
- GROUP, 27
- IncrPosnFT, 18
- LABELED_PATHNAME, 31
- legal_pathname, 9
- LegalPathname, 10
- LegalPathnames, 9
- link
 - DACops constant, 28
- LinkBAS, 24
- LinkFS, 12
- LinkOK_SFS, 40
- LinkOK_UFS, 34
- LinkPermittedDAC, 34
- LinkPermittedMAC, 40
- lpnpathname, 31
- lpntype, 31
- MAC, 37
- MAC_File_Attributes, 36
- MAC_PERM, 36
- MAC_PERMS, 36
- MAC_Process_Attributes, 37
- MACops, 37
- MACVisiblePathname, 38
- MACVisiblePathnames, 38
- MaybeTruncateFile, 17
- NewFid, 5
- NonExistentFile, 10
- OID, 14
- OidToFid, 15
- OidToPosn, 15
- OnlyOpenDirectoriesRDONLY, 24
- open
 - DACops constant, 28
- OpenBAS, 24
- OpenedForAppend, 19
- OpenedForRead, 18
- OpenedForWrite, 19
- OpenExistingFileOK_SFS, 40
- OpenExistingFileOK_UFS, 34
- OpenFT, 17
- OpenInternalFT, 16
- OpenPermittedDAC, 34
- OpenPermittedMAC, 40
- padfile, 4
- ParseDir, 7
- PATHNAME, 7
- pathname_fid, 9
- PathnameLookup, 10
- PATHNAMETYPE, 31
- PERM, 27
- PID, 19
- ProcessTable, 19
- r
 - PERM constant, 27
- rdonly
 - ACCESS_MODE constant, 14
- rdwr
 - ACCESS_MODE constant, 14

ReadBAS, 25
 ReadBFS, 6
 ReadOK_SFS, 41
 ReadOK_UFS, 35
 regular
 FTYPE constant, 7
 RegularFile, 11
 relative
 PATHNAMETYPE constant, 31
 rooted_pathname_fid, 23
 RootedPathnames, 23
 Rootfid_is_preserved, 11

 schemas, *see also* types
 AbsolutePathname, 32
 AddFType, 13
 AddProcessFD, 20
 AppendBAS, 26
 BasicAccessSystem, 23
 BasicFileSpace, 4
 CloseBAS, 24
 CloseFT, 17
 CloseOK_SFS, 41
 CloseOK_UFS, 35
 CreateBAS, 25
 CreateBFS, 5
 CreateFileAttrsDAC, 33
 CreateFileAttrsMAC, 39
 CreateFS, 14
 CreateOK_SFS, 40
 CreateOK_UFS, 33
 CreatePermittedDAC, 32
 CreatePermittedMAC, 39
 DAC, 28
 DAC_File_Attributes, 28
 DAC_Process_Attributes, 28
 DACVisiblePathname, 30
 DACVisiblePathnames, 30
 DataLength, 25
 DeleteFType, 14
 DeleteProcessFD, 20
 DestroyBFS, 5
 DestroyFileAttrDAC, 29
 DirectoryFile, 11
 DuplicateFD, 21
 EmptyDirectory, 11
 ExistentFile, 10
 ExistsFileTableReferenceToFid, 22
 ExistsLinkReferenceToFid, 22
 ExistsReferenceToFid, 22
 ExistsReferenceToFileTableEntry, 22
 FdToFid, 19
 FdToOid, 19
 FileLength, 5
 FileSpace, 8
 FileSpaceGC, 22
 FileTable, 15
 FileTableGC, 22
 ForkBAS, 27
 ForkPT, 21
 FPERMS, 27
 GetContext, 42
 GetContextOK_SFS, 42
 GetContextPermittedMAC, 42
 IncrPosnFT, 18
 LegalPathname, 10
 LegalPathnames, 9
 LinkBAS, 24
 LinkFS, 12
 LinkOK_SFS, 40
 LinkOK_UFS, 34
 LinkPermittedDAC, 34
 LinkPermittedMAC, 40
 MAC, 37
 MAC_File_Attributes, 36
 MAC_Process_Attributes, 37
 MACVisiblePathname, 38
 MACVisiblePathnames, 38
 MaybeTruncateFile, 17
 NewFid, 5
 NonExistentFile, 10
 OidToFid, 15
 OidToPosn, 15
 OnlyOpenDirectoriesRDONLY, 24
 OpenBAS, 24
 OpenedForAppend, 19
 OpenedForRead, 18
 OpenedForWrite, 19
 OpenExistingFileOK_SFS, 40
 OpenExistingFileOK_UFS, 34
 OpenFT, 17
 OpenInternalFT, 16
 OpenPermittedDAC, 34
 OpenPermittedMAC, 40
 PathnameLookup, 10
 ProcessTable, 19
 ReadBAS, 25
 ReadBFS, 6
 ReadOK_SFS, 41
 ReadOK_UFS, 35

- RegularFile, 11
- RootedPathnames, 23
- Rootfid_is_preserved, 11
- SeekEOF, 18
- SeekFT, 18
- TruncateFile, 17
- TruncatingOpen, 17
- UFS, 31, 39
- UnlinkBAS, 25
- UnlinkFS, 13
- UnlinkOK_SFS, 41
- UnlinkOK_UFS, 35
- UnlinkPermittedDAC, 35
- UnlinkPermittedMAC, 41
- WriteBAS, 26
- WriteBFS, 6
- WriteOK_SFS, 41
- WriteOK_UFS, 36
- search
 - DACops constant, 28
- SECCONTEXT, 36
- SeekEOF, 18
- SeekFT, 18
- state variables
 - cpid, 31, 39
 - fdac, 28
 - fmac, 37
 - ftfid, 15
 - ftmode, 15
 - ftposn, 15
 - pcwd, 31, 39
 - pdac, 28
 - pfhtable, 19
 - pmac, 37
- SYLLABLE, 7
- translations, *see* data translations
- TruncatingOpen, 17
- types, *see also* free types
 - ACCESS_MODE, 14
 - DAC_File_Attributes, 28
 - DAC_Process_Attributes, 28
 - DACops, 28
 - EXT_MODE, 37
 - FPERMS, 27
 - LABELED_PATHNAME, 31
 - MAC_PERM, 36
 - MAC_PERMS, 36
 - MACops, 37
 - PATHNAMETYPE, 31
 - PERM, 27
 - PERMS, 27
 - UFS, 31, 39
 - unlink
 - DACops constant, 28
 - UnlinkBAS, 25
 - UnlinkFS, 13
 - UnlinkOK_SFS, 41
 - UnlinkOK_UFS, 35
 - UnlinkPermittedDAC, 35
 - UnlinkPermittedMAC, 41
 - UnParseDir, 7
 - USER, 27
- w
 - PERM constant, 27
- WriteBAS, 26
- WriteBFS, 6
- WriteOK_SFS, 41
- WriteOK_UFS, 36
- wronly
 - ACCESS_MODE constant, 14
- x
 - PERM constant, 27