

# **Scsh Reference Manual**

---

For scsh release 0.6.5  
November 2003

**Olin Shivers, Brian D. Carlstrom, Martin Gasbichler, and Mike Sperber**

---

# Acknowledgements

Who should I thank? My so-called “colleagues,” who laugh at me behind my back, all the while becoming famous on *my* work? My worthless graduate students, whose computer skills appear to be limited to downloading bitmaps off of netnews? My parents, who are still waiting for me to quit “fooling around with computers,” go to med school, and become a radiologist? My department chairman, a manager who gives one new insight into and sympathy for disgruntled postal workers?

My God, no one could blame me—no one!—if I went off the edge and just lost it completely one day. I couldn’t get through the day as it is without the Prozac and Jack Daniels I keep on the shelf, behind my Tops-20 JSYS manuals. I start getting the shakes real bad around 10am, right before my advisor meetings. A 10 oz. Jack ’n Zac helps me get through the meetings without one of my students winding up with his severed head in a bowling-ball bag. They look at me funny; they think I twitch a lot. I’m not twitching. I’m controlling my impulse to snag my 9mm Sig-Sauer out from my day-pack and make a few strong points about the quality of undergraduate education in Amerika.

If I thought anyone cared, if I thought anyone would even be reading this, I’d probably make an effort to keep up appearances until the last possible moment. But no one does, and no one will. So I can pretty much say exactly what I think.

Oh, yes, the *acknowledgements*. I think not. I did it. I did it all, by myself.

Olin Shivers  
Cambridge  
September 4, 1994

# Contents

<b>Contents</b>	<b>iii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Copyright & source-code license . . . . .	1
1.2 Obtaining scsh . . . . .	2
1.3 Building scsh . . . . .	2
1.4 Caveats . . . . .	3
1.5 Naming conventions . . . . .	3
1.6 Lexical issues . . . . .	4
1.6.1 Extended symbol syntax . . . . .	4
1.6.2 Extended string syntax . . . . .	5
1.6.3 Block comments and executable interpreter-triggers . . .	5
1.6.4 Here-strings . . . . .	5
1.6.5 Dot . . . . .	7
1.7 Record types and the <code>define-record</code> form . . . . .	7
1.8 A word about Unix standards . . . . .	9
<b>2 Process notation</b>	<b>10</b>
2.1 Extended process forms and I/O redirections . . . . .	10
2.1.1 Port and file descriptor sync . . . . .	11
2.2 Process forms . . . . .	12
2.3 Using extended process forms in Scheme . . . . .	13
2.3.1 Procedures and special forms . . . . .	14
2.3.2 Interfacing process output to Scheme . . . . .	14
2.4 More complex process operations . . . . .	16
2.4.1 Pids and ports together . . . . .	17

2.4.2	Multiple stream capture . . . . .	17
2.5	Conditional process sequencing forms . . . . .	19
2.6	Process filters . . . . .	20
<b>3</b>	<b>System Calls</b>	<b>21</b>
3.1	Errors . . . . .	21
3.1.1	Interactive mode and error handling . . . . .	23
3.2	I/O . . . . .	24
3.2.1	Standard R5RS I/O procedures . . . . .	24
3.2.2	Port manipulation and standard ports . . . . .	24
3.2.3	String ports . . . . .	26
3.2.4	Revealed ports and file descriptors . . . . .	26
3.2.5	Port-mapping machinery . . . . .	29
3.2.6	Unix I/O . . . . .	31
3.2.7	Buffered I/O . . . . .	38
3.2.8	File locking . . . . .	39
3.3	File system . . . . .	41
3.4	Processes . . . . .	54
3.4.1	Process objects and process reaping . . . . .	57
3.4.2	Process waiting . . . . .	60
3.4.3	Analysing process status codes . . . . .	61
3.5	Process state . . . . .	62
3.6	User and group database access . . . . .	64
3.7	Accessing command-line arguments . . . . .	64
3.8	System parameters . . . . .	66
3.9	Signal system . . . . .	67
3.10	Time . . . . .	71
3.10.1	Terminology . . . . .	71
3.10.2	Basic data types . . . . .	71
3.10.3	Time zones . . . . .	73
3.10.4	Procedures . . . . .	73
3.11	Environment variables . . . . .	77
3.11.1	Path lists and colon lists . . . . .	79
3.11.2	\$USER, \$HOME, and \$PATH . . . . .	80
3.12	Terminal device control . . . . .	81

3.12.1	Portability across OS variants . . . . .	81
3.12.2	Miscellaneous procedures . . . . .	82
3.12.3	The tty-info record type . . . . .	82
3.12.4	Using tty-info records . . . . .	84
3.12.5	Other terminal-device procedures . . . . .	85
3.12.6	Control terminals, sessions, and terminal process groups	86
3.12.7	Pseudo-terminals . . . . .	87
<b>4</b>	<b>Networking</b>	<b>94</b>
4.1	High-level interface . . . . .	94
4.2	Sockets . . . . .	95
4.3	Socket addresses . . . . .	97
4.4	Socket primitives . . . . .	98
4.5	Performing input and output on sockets . . . . .	99
4.6	Socket options . . . . .	100
4.7	Database-information entries . . . . .	101
<b>5</b>	<b>Strings and characters</b>	<b>103</b>
5.1	Manipulating file names . . . . .	104
5.1.1	Terminology . . . . .	104
5.1.2	Procedures . . . . .	105
5.2	Other string manipulation facilities . . . . .	110
5.3	ASCII encoding . . . . .	110
5.4	Character predicates . . . . .	110
5.5	Deprecated character-set procedures . . . . .	111
<b>6</b>	<b>Pattern-matching strings with regular expressions</b>	<b>112</b>
6.1	Summary SRE syntax . . . . .	113
6.2	Examples . . . . .	117
6.3	A short tutorial . . . . .	117
6.4	Choices . . . . .	119
6.4.1	Embedding regexps within Scheme programs . . . . .	128
6.5	Regexp functions . . . . .	129
6.5.1	Obsolete, deprecated procedures . . . . .	129
6.5.2	Standard procedures and syntax . . . . .	129
6.6	The regexp ADT . . . . .	135

6.7	Syntax-hacking tools . . . . .	138
<b>7</b>	<b>Reading delimited strings</b>	<b>140</b>
<b>8</b>	<b>Awk, record I/O, and field parsing</b>	<b>143</b>
8.1	Record I/O and field parsing . . . . .	143
8.1.1	Reading records . . . . .	144
8.1.2	Parsing fields . . . . .	144
8.1.3	Field readers . . . . .	148
8.1.4	Forward-progress guarantees and empty-string matches	148
8.1.5	Reader limitations . . . . .	150
8.2	Awk . . . . .	150
8.2.1	Examples . . . . .	153
8.3	Backwards compatibility . . . . .	155
<b>9</b>	<b>Concurrent system programming</b>	<b>156</b>
9.1	Threads . . . . .	156
9.2	Locks . . . . .	157
9.3	Placeholders . . . . .	158
9.4	The event interface to interrupts . . . . .	158
9.5	Interaction between threads and process state . . . . .	159
<b>10</b>	<b>Miscellaneous routines</b>	<b>161</b>
10.1	Integer bitwise ops . . . . .	161
10.2	Password encryption . . . . .	161
10.3	Dot-Locking . . . . .	161
10.4	Syslog facility . . . . .	163
10.5	MD5 interface . . . . .	167
<b>11</b>	<b>Running scsh</b>	<b>168</b>
11.1	Scsh command-line switches . . . . .	169
11.1.1	Scripts and programs . . . . .	169
11.1.2	Inserting interpreter triggers into scsh programs . . . . .	169
11.1.3	Module system . . . . .	170
11.1.4	Switches . . . . .	171
11.1.5	The meta argument . . . . .	177
11.1.6	Examples . . . . .	179

11.1.7	Process exit values . . . . .	181
11.2	The scsh virtual machine . . . . .	182
11.2.1	VM arguments . . . . .	182
11.2.2	Stripped image . . . . .	183
11.2.3	Inserting interpreter triggers into heap images . . . . .	184
11.2.4	Inserting a double-level trigger into Scheme programs . . . . .	184
11.3	Compiling scsh programs . . . . .	184
11.4	Standard file locations . . . . .	185
	<b>Index</b>	<b>187</b>





# Chapter 1

## Introduction

This is the reference manual for `scsh`, a Unix shell that is embedded within Scheme. `Scsh` is a Scheme system designed for writing useful standalone Unix programs and shell scripts—it spans a wide range of application, from “script” applications usually handled with `perl` or `sh`, to more standard systems applications usually written in C.

`Scsh` comes built on top of Scheme 48, and has two components: a process notation for running programs and setting up pipelines and redirections, and a complete syscall library for low-level access to the operating system. This manual gives a complete description of `scsh`. A general discussion of the design principles behind `scsh` can be found in a companion paper, “A Scheme Shell.”

### 1.1 Copyright & source-code license

`Scsh` is open source. The complete sources come with the standard distribution, which can be downloaded off the net. `Scsh` has an ideologically hip, BSD-style license.

We note that the code is a rich source for other Scheme implementations to mine. Not only the *code*, but the *APIs* are available for implementors working on Scheme environments for systems programming. These APIs represent years of work, and should provide a big head-start on any related effort. (Just don’t call it “`scsh`,” unless it’s *exactly* compliant with the `scsh` interfaces.)

Take all the code you like; we’ll just write more.

## 1.2 Obtaining scsh

Scsh is distributed via net publication. We place new releases at well-known network sites, and allow them to propagate from there. We currently release scsh to the following Internet sites:

```
ftp://ftp.scsh.net/pub/scsh
http://prdownloads.sourceforge.net/scsh/
```

Each should have a compressed tar file of the entire scsh release, which includes all the source code and the manual, and a separate file containing just this manual in Postscript form, for those who simply wish to read about the system.

However, nothing is certain for long on the Net. Probably the best way to get a copy of scsh is to use a network resource-discovery tool, such as archie, to find ftp servers storing scsh tar files. Take the set of sites storing the most recent release of scsh, choose one close to your site, and download the tar file.

## 1.3 Building scsh

Scsh currently runs on a fairly large set of Unix systems, including Linux, FreeBSD, OpenBSD, NetBSD, MacOS X, SunOS, Solaris, AIX, NeXTSTEP, Irix, and HP-UX. We use the Gnu project's autoconfig tool to generate self-configuring shell scripts that customise the scsh Makefile for different OS variants. This means that if you use one of the common Unix implementations, building scsh should require exactly the following steps:

<code>gunzip scsh.tar.gz</code>	<i>Uncompress the release tar file.</i>
<code>untar xfv scsh.tar</code>	<i>Unpack the source code.</i>
<code>cd scsh-0.6.x</code>	<i>Move to the source directory.</i>
<code>./configure</code>	<i>Examine host; build Makefile.</i>
<code>make</code>	<i>Build system.</i>

When you are done, you should have a virtual machine compiled in file `scshvm`, and a heap image in file `scsh/scsh.image`. Typing

```
make install
```

will install these programs in your installation directory (by default, `/usr/local`), along with a small stub startup binary, `scsh`.

If you don't have the patience to do this, you can start up a Scheme shell immediately after the initial make by simply saying

```
./scshvm -o ./scshvm -i scsh/scsh.image
```

See chapter 11 for full details on installation locations and startup options.

It is not too difficult to port `scsh` to another Unix platform if your OS is not supported by the current release. See the release notes for more details on how to do this.

## 1.4 Caveats

It is important to note what `scsh` is *not*, as well as what it is. `Scsh`, in the current release, is primarily designed for the writing of shell scripts—programming. It is not a very comfortable system for interactive command use: the current release lacks job control, command-line editing, a terse, convenient command syntax, and it does not read in an initialisation file analogous to `.login` or `.profile`. We hope to address all of these issues in future releases; we even have designs for several of these features; but the system as-released does not currently provide these features.

## 1.5 Naming conventions

`Scsh` follows a general naming scheme that consistently employs a set of abbreviations. This is intended to make it easier to remember the names of things. Some of the common ones are:

`fdes` Means “file descriptor,” a small integer used in Unix to represent I/O channels.

...\* A given bit of functionality sometimes comes in two related forms, the first being a *special form* that contains a body of Scheme code to be executed in some context, and the other being a *procedure* that takes a procedural argument (a “thunk”) to be called in the same context. The procedure variant is named by taking the name of the special form, and appending an asterisk. For example:

```
;;; Special form:
(with-cwd "/etc"
  (for-each print-file (directory-files))
  (display "All done"))

;;; Procedure:
(with-cwd* "/etc"
  (lambda ()
    (for-each print-file (directory-files))
    (display "All done")))
```

*action/modifier* The infix “/” is pronounced “with,” as in `exec/env`—“exec with environment.”

`call/...` Procedures that call their argument on some computed value are usually named “`call/...`,” e.g., `(call/fdes port proc)`, which calls `proc` on `port`’s file descriptor, returning whatever `proc` returns. The abbreviated name means “call with file descriptor.”

`with-...` Procedures that call their argument, and special forms that execute their bodies in some special dynamic context frequently have names of the form `with-....`. For example, `(with-env env body1 ...)` and `(with-env* env thunk)`. These forms set the process environment body, execute their body or thunk, and then return after resetting the environment to its original state.

`create-` Procedures that create objects in the file system (files, directories, temp files, fifos, etc.), begin with `create-....`

`delete-` Procedures that delete objects from the file system (files, directories, temp files, fifos, etc.), begin with `delete-....`

*record:field* Procedures that access fields of a record are usually written with a colon between the name of the record and the name of the field, as in `user-info:home-dir`.

`%...` A percent sign is used to prefix lower-level scsh primitives that are not commonly used.

`-info` Data structures packaging up information about various OS entities frequently end in `...-info`. Examples: `user-info`, `file-info`, `group-info`, and `host-info`.

Enumerated constants from some set *s* are usually named `s/const1`, `s/const2`, .... For example, the various Unix signal integers have the names `signal/cont`, `signal/kill`, `signal/int`, `signal/hup`, and so forth.

## 1.6 Lexical issues

Scsh’s lexical syntax is just R5RS Scheme, with the following exceptions.

### 1.6.1 Extended symbol syntax

Scsh’s symbol syntax differs from R5RS Scheme in the following ways:

- In scsh, symbol case is preserved by `read` and is significant on symbol comparison. This means

```
(run (less Readme))
```

displays the right file.

- “-” and “+” are allowed to begin symbols. So the following are legitimate symbols:

```
-02 -geometry +Wn
```

- “|” and “.” are symbol constituents. This allows | for the pipe symbol, and .. for the parent-directory symbol. (Of course, “.” alone is not a symbol, but a dotted-pair marker.)
- A symbol may begin with a digit. So the following are legitimate symbols:

```
9x15 80x36-3+440
```

### 1.6.2 Extended string syntax

Scsh strings are allowed to contain the ANSI C escape sequences such as `\n` and `\161`.

### 1.6.3 Block comments and executable interpreter-triggers

Scsh allows source files to begin with a header of the form

```
#!/usr/local/bin/scsh -s
```

The Unix operating system treats source files beginning with the headers of this form specially; they can be directly executed by the operating system (see chapter 11 for information on how to use this feature). The scsh interpreter ignores this special header by treating `#!` as a comment marker similar to `;`. When the scsh reader encounters `#!`, it skips characters until it finds the closing sequence `newline/exclamation-point/sharp-sign/newline`.

Although the form of the `#!` read-macro was chosen to support interpreter-triggers for executable Unix scripts, it is a general block-comment sequence and can be used as such anywhere in a scsh program.

### 1.6.4 Here-strings

The read macro `#<` is used to introduce “here-strings” in programs, similar to the `<<` “here document” redirections provided by `sh` and `csh`. There are two kinds of here-string, character-delimited and line-delimited; they are both introduced by the `#<` sequence.

## Character-delimited here-strings

A *character-delimited* here-string has the form

```
#<x...stuff...x
```

where *x* is any single character (except <, see below), which is used to delimit the string bounds. Some examples:

Here-string syntax	Ordinary string syntax
#< Hello, world.	"Hello, world."
#<!"Ouch," he said.!	"\"Ouch,\" he said."

There is no interpretation of characters within the here-string; the characters are all copied verbatim.

## Line-delimited here-strings

If the sequence begins "#<<" then it introduces a *line-delimited* here-string. These are similar to the "here documents" of sh and csh. Line-delimited here-strings are delimited by the rest of the text line that follows the "#<<" sequence. For example:

```
#<<FOO
Hello, there.
This is read by Scheme as a string,
terminated by the first occurrence
of newline-F-0-0-newline or newline-F-0-0-eof.
FOO
```

Thus,

```
#<<foo
Hello, world.
foo
```

is the same thing as

```
"Hello, world."
```

Line-delimited here-strings are useful for writing down long, constant strings—such as long, multi-line format strings, or arguments to Unix programs, *e.g.*,

```
;; Free up some disk space for my netnews files.
(run (csh -c #<<EOF
cd /uops
rm -rf *
echo All done.

EOF
))
```

The advantage they have over the double-quote syntax (e.g., "Hello, world.") is that there is no need to backslash-quote special characters internal to the string, such as the double-quote or backslash characters.

The detailed syntax of line-delimited here-strings is as follows. The characters "#<<" begin the here-string. The characters between the "#<<" and the next newline are the *delimiter line*. All characters between the "#<<" and the next newline comprise the delimiter line—including any white space. The body of the string begins on the following line, and is terminated by a line of text which exactly matches the delimiter line. This terminating line can be ended by either a newline or end-of-file. Absolutely no interpretation is done on the input string. Control characters, white space, quotes, backslash—everything is copied as-is. The newline immediately preceding the terminating delimiter line is not included in the result string (leave an extra blank line if you need to put a final newline in the here-string—see the example above). If EOF is encountered before reading the end of the here-string, an error is signalled.

### 1.6.5 Dot

It is unfortunate that the single-dot token, ".", is both a fundamental Unix file name and a deep, primitive syntactic token in Scheme—it means the following will not parse correctly in scsh:

```
(run/strings (find . -name *.c -print))
```

You must instead quote the dot:

```
(run/strings (find "." -name *.c -print))
```

When you write shell scripts that manipulate the file system, keep in mind the special status of the dot token.

## 1.7 Record types and the define-record form

Scsh's interfaces occasionally provide data in structured record types; an example is the `file-info` record whose various fields describe the size, protection,

last date of modification, and other pertinent data for a particular file. These record types are described in this manual using the `define-record` notation, which looks like the following:

```
(define-record ship
  x
  y
  (size 100))
```

This form defines a *ship* record, with three fields: its *x* and *y* coordinates, and its size. The values of the *x* and *y* fields are specified as parameters to the ship-building procedure, `(make-ship x y)`, and the *size* field is initialised to 100. All told, the `define-record` form above defines the following procedures:

Procedure	Definition
<code>(make-ship x y)</code>	Create a new <i>ship</i> record.
<code>(ship:x ship)</code>	Retrieve the <i>x</i> field.
<code>(ship:y ship)</code>	Retrieve the <i>y</i> field.
<code>(ship:size ship)</code>	Retrieve the <i>size</i> field.
<code>(set-ship:x ship new-x)</code>	Assign the <i>x</i> field.
<code>(set-ship:y ship new-y)</code>	Assign the <i>y</i> field.
<code>(set-ship:size ship new-size)</code>	Assign the <i>size</i> field.
<code>(modify-ship:x ship xfun)</code>	Modify <i>x</i> field with <i>xfun</i> .
<code>(modify-ship:y ship yfun)</code>	Modify <i>y</i> field with <i>yfun</i> .
<code>(modify-ship:size ship sizefun)</code>	Modify <i>size</i> field with <i>sizefun</i> .
<code>(ship? object)</code>	Type predicate.
<code>(copy-ship ship)</code>	Shallow-copy of the record.

An implementation of `define-record` is available as a macro for Scheme programmers to define their own record types; the syntax is accessed by opening the package `defrec-package`, which exports the single syntax form `define-record`. See the source code for the `defrec-package` module for further details of the macro.

You must open this package to access the form. `Scsh` does not export a record-definition package by default as there are several from which to choose. Besides the `define-record` macro, which Shivers prefers<sup>1</sup>, you might instead wish to employ the notationally-distinct `define-record-type` macro that Jonathan Rees prefers<sup>2</sup>. It can be found in the `define-record-types` structure.

Alternatively, you may define your own, of course.

---

<sup>1</sup>He wrote it.

<sup>2</sup>He wrote it.



## 1.8 A word about Unix standards

“The wonderful thing about Unix standards is that there are so many to choose from.” You may be totally bewildered about the multitude of various standards that exist. Rest assured that nowhere in this manual will you encounter an attempt to spell it all out for you; you could not read and internalise such a twisted account without bleeding from the nose and ears.

However, you might keep in mind the following simple fact: of all the standards, POSIX is the least common denominator. So when this manual repeatedly refers to POSIX, the point is “the thing we are describing should be portable just about anywhere.” Scsh sticks to POSIX when at all possible; its major departure is symbolic links, which aren’t in POSIX (see—it really *is* a least common denominator).

## Chapter 2

# Process notation

Scsh has a notation for controlling Unix processes that takes the form of s-expressions; this notation can then be embedded inside of standard Scheme code. The basic elements of this notation are *process forms*, *extended process forms*, and *redirections*.

### 2.1 Extended process forms and I/O redirections

An *extended process form* is a specification of a Unix process to run, in a particular I/O environment:

$$epf ::= (pf \textit{redir}_1 \dots \textit{redir}_n)$$

where *pf* is a process form and the *redir<sub>i</sub>* are redirection specs. A *redirection spec* is one of:

(< [ <i>fdes</i> ] <i>file-name</i> )	; Open file for read.
(> [ <i>fdes</i> ] <i>file-name</i> )	; Open file create/truncate.
(<< [ <i>fdes</i> ] <i>object</i> )	; Use <i>object</i> 's printed rep.
(>> [ <i>fdes</i> ] <i>file-name</i> )	; Open file for append.
(= <i>fdes</i> <i>fdes/port</i> )	; Dup2
(- <i>fdes/port</i> )	; Close <i>fdes/port</i> .
stdports	; 0,1,2 dup'd from standard ports.

The input redirections default to file descriptor 0; the output redirections default to file descriptor 1.

The subforms of a redirection are implicitly backquoted, and symbols stand for their print-names. So (> ,x) means "output to the file named by Scheme variable x," and (< /usr/shivers/.login) means "read from /usr/shivers/.login."

Here are two more examples of I/O redirection:

```
(< ,(vector-ref fv i))
(>> 2 /tmp/buf)
```

These two redirections cause the file `fv[i]` to be opened on `stdin`, and `/tmp/buf` to be opened for append writes on `stderr`.

The redirection `(<< object)` causes input to come from the printed representation of *object*. For example,

```
(<< "The quick brown fox jumped over the lazy dog.")
```

causes reads from `stdin` to produce the characters of the above string. The object is converted to its printed representation using the `display` procedure, so

```
(<< (A five element list))
```

is the same as

```
(<< "(A five element list)")
```

is the same as

```
(<< ,(reverse '(list element five A))).
```

(Here we use the implicit backquoting feature to compute the list to be printed.)

The redirection `(= fdes fdes/port)` causes *fdes/port* to be dup'd into file descriptor *fdes*. For example, the redirection

```
(= 2 1)
```

causes `stderr` to be the same as `stdout`. *fdes/port* can also be a port, for example:

```
(= 2 ,(current-output-port))
```

causes `stderr` to be dup'd from the current output port. In this case, it is an error if the port is not a file port (*e.g.*, a string port). More complex redirections can be accomplished using the `begin process` form, discussed below, which gives the programmer full control of I/O redirection from Scheme.

### 2.1.1 Port and file descriptor sync

It's important to remember that rebinding Scheme's current I/O ports (*e.g.*, using `call-with-input-file` to rebind the value of `(current-input-port)`) does *not* automatically "rebind" the file referenced by the Unix `stdio` file descriptors 0, 1, and 2. This is impossible to do in general, since some Scheme ports are not representable as Unix file descriptors. For example, many Scheme implementations provide "string ports," that is, ports that collect characters sent to them into memory buffers. The accumulated string

can later be retrieved from the port as a string. If a user were to bind (current-output-port) to such a port, it would be impossible to associate file descriptor 1 with this port, as it cannot be represented in Unix. So, if the user subsequently forked off some other program as a subprocess, that program would of course not see the Scheme string port as its standard output.

To keep stdio synced with the values of Scheme's current I/O ports, use the special redirection `stdports`. This causes 0, 1, 2 to be redirected from the current Scheme standard ports. It is equivalent to the three redirections:

```
(= 0 ,(current-input-port))
(= 1 ,(current-output-port))
(= 2 ,(error-output-port))
```

The redirections are done in the indicated order. This will cause an error if one of the current I/O ports isn't a Unix port (e.g., if one is a string port). This Scheme/Unix I/O synchronisation can also be had in Scheme code (as opposed to a redirection spec) with the `(stdports->stdio)` procedure.

## 2.2 Process forms

A *process form* specifies a computation to perform as an independent Unix process. It can be one of the following:

```
(begin . scheme-code)           ; Run scheme-code in a fork.
(| pf1 ... pfn)                 ; Simple pipeline
(|+ connect-list pf1 ... pfn) ; Complex pipeline
(epf . epf)                     ; An extended process form.
(prog arg1 ... argn)           ; Default: exec the program.
```

The default case `(prog arg1 ... argn)` is also implicitly backquoted. That is, it is equivalent to:

```
(begin (apply exec-path '(prog arg1 ... argn)))
```

Exec-path is the version of the `exec()` system call that uses `scsh`'s path list to search for an executable. The program and the arguments must be either strings, symbols, or integers. Symbols and integers are coerced to strings. A symbol's print-name is used. Integers are converted to strings in base 10. Using symbols instead of strings is convenient, since it suppresses the clutter of the surrounding "... " quotation marks. To aid this purpose, `scsh` reads symbols in a case-sensitive manner, so that you can say

```
(more Readme)
```

and get the right file.

A *connect-list* is a specification of how two processes are to be wired together by pipes. It has the form `((from1 from2 ... to) ...)` and is implicitly backquoted. For example,

```
(|+ ((1 2 0) (3 1)) pf1 pf2)
```

runs *pf<sub>1</sub>* and *pf<sub>2</sub>*. The first clause `(1 2 0)` causes *pf<sub>1</sub>*'s stdout (1) and stderr (2) to be connected via pipe to *pf<sub>2</sub>*'s stdin (0). The second clause `(3 1)` causes *pf<sub>1</sub>*'s file descriptor 3 to be connected to *pf<sub>2</sub>*'s file descriptor 1.

The *begin* process form does a `stdio->stdports` synchronisation in the child process before executing the body of the form. This guarantees that the *begin* form, like all other process forms, “sees” the effects of any associated I/O redirections.

Note that R5RS does not specify whether or not `|` and `|+` are readable symbols. Scsh does.

## 2.3 Using extended process forms in Scheme

Process forms and extended process forms are *not* Scheme. They are a different notation for expressing computation that, like Scheme, is based upon s-expressions. Extended process forms are used in Scheme programs by embedding them inside special Scheme forms. There are three basic Scheme forms that use extended process forms: `exec-epf`, `&`, and `run`.

<code>(exec-epf . epf)</code>	$\longrightarrow$	<i>no return value</i>	syntax
<code>(&amp; . epf)</code>	$\longrightarrow$	<i>proc</i>	syntax
<code>(run . epf)</code>	$\longrightarrow$	<i>status</i>	syntax

The `(exec-epf . epf)` form nukes the current process: it establishes the I/O redirections and then overlays the current process with the requested computation.

The `(& . epf)` form is similar, except that the process is forked off in background. The form returns the subprocess' process object.

The `(run . epf)` form runs the process in foreground: after forking off the computation, it waits for the subprocess to exit, and returns its exit status.

These special forms are macros that expand into the equivalent series of system calls. The definition of the `exec-epf` macro is non-trivial, as it produces the code to handle I/O redirections and set up pipelines. However, the definitions of the `&` and `run` macros are very simple:

```
(& . epf)  ≡ (fork (λ () (exec-epf . epf)))
(run . epf) ≡ (wait (& . epf))
```

### 2.3.1 Procedures and special forms

It is a general design principle in *scsh* that all functionality made available through special syntax is also available in a straightforward procedural form. So there are procedural equivalents for all of the process notation. In this way, the programmer is not restricted by the particular details of the syntax. Here are some of the syntax/procedure equivalents:

Notation	Procedure
	fork/pipe
+	fork/pipe+
exec-epf	exec-path
redirection	open, dup
&	fork
run	wait + fork

Having a solid procedural foundation also allows for general notational experimentation using Scheme's macros. For example, the programmer can build his own pipeline notation on top of the `fork` and `fork/pipe` procedures. Chapter 3 gives the full story on all the procedures in the `syscall` library.

### 2.3.2 Interfacing process output to Scheme

There is a family of procedures and special forms that can be used to capture the output of processes as Scheme data.

(run/port .epf)	→	port	syntax
(run/file .epf)	→	string	syntax
(run/string .epf)	→	string	syntax
(run/strings .epf)	→	string list	syntax
(run/sexp .epf)	→	object	syntax
(run/sexps .epf)	→	list	syntax

These forms all fork off subprocesses, collecting the process' output to stdout in some form or another. The subprocess runs with file descriptor 1 and the current output port bound to a pipe.

<code>run/port</code>	Value is a port open on process's stdout. Returns immediately after forking child.
<code>run/file</code>	Value is name of a temp file containing process's output. Returns when process exits.
<code>run/string</code>	Value is a string containing process' output. Returns when eof read.
<code>run/strings</code>	Splits process' output into a list of newline-delimited strings. Returns when eof read.
<code>run/sexp</code>	Reads a single object from process' stdout with read. Returns as soon as the read completes.
<code>run/sexps</code>	Repeatedly reads objects from process' stdout with read. Returns accumulated list upon eof.

The delimiting newlines are not included in the strings returned by `run/strings`.

These special forms just expand into calls to the following analogous procedures.

<code>(run/port* thunk)</code>	$\longrightarrow$ <i>port</i>	procedure
<code>(run/file* thunk)</code>	$\longrightarrow$ <i>string</i>	procedure
<code>(run/string* thunk)</code>	$\longrightarrow$ <i>string</i>	procedure
<code>(run/strings* thunk)</code>	$\longrightarrow$ <i>string list</i>	procedure
<code>(run/sexp* thunk)</code>	$\longrightarrow$ <i>object</i>	procedure
<code>(run/sexps* thunk)</code>	$\longrightarrow$ <i>object list</i>	procedure

For example, `(run/port . epf)` expands into

`(run/port* (λ () (exec-epf . epf)))`.

The following procedures are also of utility for generally parsing input streams in scsh:

<code>(port-&gt;string port)</code>	$\longrightarrow$ <i>string</i>	procedure
<code>(port-&gt;sexp-list port)</code>	$\longrightarrow$ <i>list</i>	procedure
<code>(port-&gt;string-list port)</code>	$\longrightarrow$ <i>string list</i>	procedure
<code>(port-&gt;list reader port)</code>	$\longrightarrow$ <i>list</i>	procedure

`Port->string` reads the port until eof, then returns the accumulated string. `Port->sexp-list` repeatedly reads data from the port until eof, then returns the accumulated list of items. `Port->string-list` repeatedly reads newline-terminated strings from the port until eof, then returns the accumulated list of strings. The delimiting newlines are not

part of the returned strings. `port->list` generalises these two procedures. It uses *reader* to repeatedly read objects from a port. It accumulates these objects into a list, which is returned upon eof. The `port->string-list` and `port->sexp-list` procedures are trivial to define, being merely `port->list` curried with the appropriate parsers:

```
(port->string-list port) ≡ (port->list read-line port)
(port->sexp-list   port) ≡ (port->list read port)
```

The following compositions also hold:

```
run/string*  ≡ port->string      ◦ run/port*
run/strings* ≡ port->string-list ◦ run/port*
run/sexp*    ≡ read             ◦ run/port*
run/sexps*   ≡ port->sexp-list   ◦ run/port*
```

`(port-fold port reader op . seeds)`     $\longrightarrow$     *object\**    procedure

This procedure can be used to perform a variety of iterative operations over an input stream. It repeatedly uses *reader* to read an object from *port*. If the first read returns eof, then the entire `port-fold` operation returns the seeds as multiple values. If the first read operation returns some other value *v*, then *op* is applied to *v* and the seeds: `(op v . seeds)`. This should return a new set of seed values, and the reduction then loops, reading a new value from the port, and so forth. (If multiple seed values are used, then *op* must return multiple values.)

For example, `(port->list reader port)` could be defined as

```
(reverse (port-fold port reader cons '()))
```

An imperative way to look at `port-fold` is to say that it abstracts the idea of a loop over a stream of values read from some port, where the seed values express the loop state.

*Remark:* This procedure was formerly named `reduce-port`. The old binding is still provided, but is deprecated and will probably vanish in a future release.

## 2.4 More complex process operations

The procedures and special forms in the previous section provide for the common case, where the programmer is only interested in the output of the process. These special forms and procedures provide more complicated facilities for manipulating processes.



### 2.4.1 Pids and ports together

(run/port+proc .epf)	→	[port proc]	syntax
(run/port+proc* thunk)	→	[port proc]	procedure

This special form and its analogous procedure can be used if the programmer also wishes access to the process' pid, exit status, or other information. They both fork off a subprocess, returning two values: a port open on the process' stdout (and current output port), and the subprocess's process object. A process object encapsulates the subprocess' process id and exit code; it is the value passed to the `wait` system call.

For example, to uncompress a tech report, reading the uncompressed data into scsh, and also be able to track the exit status of the decompression process, use the following:

```
(receive (port child) (run/port+proc (zcat tr91-145.tex.Z))
  (let* ((paper (port->string port))
        (status (wait child)))
    ...use paper, status, and child here...))
```

Note that you must *first* do the `port->string` and *then* do the `wait`—the other way around may lock up when the `zcat` fills up its output pipe buffer.

### 2.4.2 Multiple stream capture

Occasionally, the programmer may want to capture multiple distinct output streams from a process. For instance, he may wish to read the stdout and stderr streams into two distinct strings. This is accomplished with the `run/collecting` form and its analogous procedure, `run/collecting*`.

(run/collecting <i>fds.epf</i> )	→	[ <i>status port...</i> ]	syntax
(run/collecting* <i>fds thunk</i> )	→	[ <i>status port...</i> ]	procedure

Run/collecting and run/collecting\* run processes that produce multiple output streams and return ports open on these streams. To avoid issues of deadlock, run/collecting doesn't use pipes. Instead, it first runs the process with output to temp files, then returns ports open on the temp files. For example,

```
(run/collecting (1 2) (ls))
```

runs `ls` with `stdout` (fd 1) and `stderr` (fd 2) redirected to temporary files. When the `ls` is done, `run/collecting` returns three values: the `ls` process' exit status, and two ports open on the temporary files. The files are deleted before `run/collecting` returns, so when the ports are closed,

they vanish. The `fds` list of file descriptors is implicitly backquoted by the special-form version.

For example, if Kaiming has his mailbox protected, then

```
(receive (status out err)
         (run/collecting (1 2) (cat /usr/kmshea/mbox))
         (list status (port->string out) (port->string err)))
```

might produce the list

```
(256 "" "cat: /usr/kmshea/mbox: Permission denied")
```

What is the deadlock hazard that causes `run/collecting` to use temp files? Processes with multiple output streams can lock up if they use pipes to communicate with Scheme I/O readers. For example, suppose some Unix program `myprog` does the following:

1. First, outputs a single "(" to `stderr`.
2. Then, outputs a megabyte of data to `stdout`.
3. Finally, outputs a single ")" to `stderr`, and exits.

Our `scsh` programmer decides to run `myprog` with `stdout` and `stderr` redirected *via Unix pipes* to the ports `port1` and `port2`, respectively. He gets into trouble when he subsequently says `(read port2)`. The Scheme `read` routine reads the open paren, and then hangs in a `read()` system call trying to read a matching close paren. But before `myprog` sends the close paren down the `stderr` pipe, it first tries to write a megabyte of data to the `stdout` pipe. However, Scheme is not reading that pipe—it's stuck waiting for input on `stderr`. So the `stdout` pipe quickly fills up, and `myprog` hangs, waiting for the pipe to drain. The `myprog` child is stuck in a `stdout/port1` write; the Scheme parent is stuck in a `stderr/port2` read. Deadlock.

Here's a concrete example that does exactly the above:

```

(receive (status port1 port2)
  (run/collecting (1 2)
    (begin
      ;; Write an open paren to stderr.
      (run (echo "(") (= 1 2))
      ;; Copy a lot of stuff to stdout.
      (run (cat /usr/dict/words))
      ;; Write a close paren to stderr.
      (run (echo ")") (= 1 2))))

;; OK. Here, I have a port PORT1 built over a pipe
;; connected to the BEGIN subproc's stdout, and
;; PORT2 built over a pipe connected to the BEGIN
;; subproc's stderr.
(read port2) ; Should return the empty list.
(port->string port1) ; Should return a big string.

```

In order to avoid this problem, `run/collecting` and `run/collecting*` first run the child process to completion, buffering all the output streams in temp files (using the `temp-file-channel` procedure, see below). When the child process exits, ports open on the buffered output are returned. This approach has two disadvantages over using pipes:

- The total output from the child output is temporarily written to the disk before returning from `run/collecting`. If this output is some large intermediate result, the disk could fill up.
- The child producer and Scheme consumer are serialised; there is no concurrency overlap in their execution.

However, it remains a simple solution that avoids deadlock. More sophisticated solutions can easily be programmed up as needed—`run/collecting*` itself is only 12 lines of simple code.

See `temp-file-channel` for more information on creating temp files as communication channels.

## 2.5 Conditional process sequencing forms

These forms allow conditional execution of a sequence of processes.

$(|| pf_1 \dots pf_n) \longrightarrow \text{boolean}$  syntax

Run each proc until one completes successfully (*i.e.*, exit status zero). Return true if some proc completes successfully; otherwise #f.

`(&& pf1 ... pfn)`  $\longrightarrow$  *boolean* syntax

Run each proc until one fails (*i.e.*, exit status non-zero). Return true if all procs complete successfully; otherwise #f.

## 2.6 Process filters

These procedures are useful for forking off processes to filter text streams.

`(make-char-port-filter filter)`  $\longrightarrow$  *procedure* procedure

The *filter* argument is a character $\rightarrow$ character procedure. Returns a procedure that when called, repeatedly reads a character from the current input port, applies *filter* to the character, and writes the result to the current output port. The procedure returns upon reaching eof on the input port.

For example, to downcase a stream of text in a spell-checking pipeline, instead of using the Unix `tr A-Z a-z` command, we can say:

```
(run (| (delatex)
      (begin ((char-filter char-downcase))) ; tr A-Z a-z
      (spell)
      (sort)
      (uniq))
    (< scsh.tex)
    (> spell-errors.txt))
```

`(make-string-port-filter filter [buflen])`  $\longrightarrow$  *procedure* procedure

The *filter* argument is a string $\rightarrow$ string procedure. Returns a procedure that when called, repeatedly reads a string from the current input port, applies *filter* to the string, and writes the result to the current output port. The procedure returns upon reaching eof on the input port.

The optional *buflen* argument controls the number of characters each internal read operation requests; this means that *filter* will never be applied to a string longer than *buflen* chars. The default *buflen* value is 1024.

## Chapter 3

# System Calls

Scsh provides (almost) complete access to the basic Unix kernel services: processes, files, signals and so forth. These procedures comprise a Scheme binding for POSIX, with a few of the more standard extensions thrown in (*e.g.*, symbolic links, `fchown`, `fstat`, sockets).

### 3.1 Errors

Scsh syscalls never return error codes, and do not use a global `errno` variable to report errors. Errors are consistently reported by raising exceptions. This frees up the procedures to return useful values, and allows the programmer to assume that *if a syscall returns, it succeeded*. This greatly simplifies the flow of the code from the programmer's point of view.

Since Scheme does not yet have a standard exception system, the scsh definition remains somewhat vague on the actual form of exceptions and exception handlers. When a standard exception system is defined, scsh will move to it. For now, scsh uses the Scheme 48 exception system, with a simple sugaring on top to hide the details in the common case.

System call error exceptions contain the Unix `errno` code reported by the system call. Unlike C, the `errno` value is a part of the exception packet, it is *not* accessed through a global variable.

For reference purposes, the Unix `errno` numbers are bound to the variables `errno/perm`, `errno/noent`, *etc.* System calls never return `error/intr`—they automatically retry.

`(errno-error errno syscall . data)`  $\longrightarrow$  *no return value* procedure

Raises a Unix error exception for Unix error number *errno*. The *syscall* and *data* arguments are packaged up in the exception packet passed to the exception handler.

(with-errno-handler\* *handler thunk*)  $\longrightarrow$  *value(s) of thunk*      procedure  
 (with-errno-handler *handler-spec . body*)  $\longrightarrow$  *value of body*      syntax

Unix syscalls raise error exceptions by calling `errno-error`. Programs can use `with-errno-handler*` to establish handlers for these exceptions.

If a Unix error arises while *thunk* is executing, *handler* is called on two arguments like this:

(*handler errno packet*)

*packet* is a list of the form

*packet* = (*errno-msg syscall . data*),

where *errno-msg* is the standard Unix error message for the error, *syscall* is the procedure that generated the error, and *data* is a list of information generated by the error, which varies from syscall to syscall.

If *handler* returns, the handler search continues upwards. *Handler* can acquire the exception by invoking a saved continuation. This procedure can be sugared over with the following syntax:

(with-errno-handler  
   ((*errno packet*) *clause ...*)  
   *body1*  
   *body2*  
   ...)

This form executes the body forms with a particular `errno` handler installed. When an `errno` error is raised, the handler search machinery will bind variable *errno* to the error's integer code, and variable *packet* to the error's auxiliary data packet. Then, the clauses will be checked for a match. The first clause that matches is executed, and its value is the value of the entire `with-errno-handler` form. If no clause matches, the handler search continues.

Error clauses have two forms

((*errno ...*) *body ...*)  
 (else *body ...*)

In the first type of clause, the *errno* forms are integer expressions. They are evaluated and compared to the error's `errno` value. An `else` clause

matches any *errno* value. Note that the *errno* and *data* variables are lexically visible to the error clauses.

Example:

```
(with-errno-handler
  ((errno packet) ; Only handle 3 particular errors.
   ((errno/wouldblock errno/again)
    (loop))
   ((errno/acces)
    (format #t "Not allowed access!")
    #f))

  (foo frobbotz)
  (blatz garglemumph))
```

It is not defined what dynamic context the handler executes in, so fluid variables cannot reliably be referenced.

Note that Scsh system calls always retry when interrupted, so that the *errno*/*intr* exception is never raised. If the programmer wishes to abort a system call on an interrupt, he should have the interrupt handler explicitly raise an exception or invoke a stored continuation to throw out of the system call.

### 3.1.1 Interactive mode and error handling

Scsh runs in two modes: interactive and script mode. It starts up in interactive mode if the scsh interpreter is started up with no script argument. Otherwise, scsh starts up in script mode. The mode determines whether scsh prints prompts in between reading and evaluating forms, and it affects the default error handler. In interactive mode, the default error handler will report the error, and generate an interactive breakpoint so that the user can interact with the system to examine, fix, or dismiss from the error. In script mode, the default error handler causes the scsh process to exit.

When scsh forks a child with (*fork*), the child resets to script mode. This can be overridden if the programmer wishes.

## 3.2 I/O

### 3.2.1 Standard R5RS I/O procedures

In `scsh`, most standard R5RS I/O operations (such as `display` or `read-char`) work on both integer file descriptors and Scheme ports. When doing I/O with a file descriptor, the I/O operation is done directly on the file, bypassing any buffered data that may have accumulated in an associated port. Note that character-at-a-time operations such as `read-char` are likely to be quite slow when performed directly upon file descriptors.

The standard R5RS procedures `read-char`, `char-ready?`, `write`, `display`, `newline`, and `write-char` are all generic, accepting integer file descriptor arguments as well as ports. `Scsh` also mandates the availability of `format`, and further requires `format` to accept file descriptor arguments as well as ports.

The procedures `peek-char` and `read` do *not* accept file descriptor arguments, since these functions require the ability to read ahead in the input stream, a feature not supported by Unix I/O.

### 3.2.2 Port manipulation and standard ports

`(close-after port consumer)`  $\longrightarrow$  *value(s) of consumer* procedure  
Returns *(consumer port)*, but closes the port on return. No dynamic-wind magic.

*Remark:* Is there a less-awkward name?

`(error-output-port)`  $\longrightarrow$  *port* procedure  
This procedure is analogous to `current-output-port`, but produces a port used for error messages—the `scsh` equivalent of `stderr`.

`(with-current-input-port* port thunk)`  $\longrightarrow$  *value(s) of thunk* procedure  
`(with-current-output-port* port thunk)`  $\longrightarrow$  *value(s) of thunk* procedure  
`(with-error-output-port* port thunk)`  $\longrightarrow$  *value(s) of thunk* procedure

These procedures install *port* as the current input, current output, and error output port, respectively, for the duration of a call to *thunk*.

`(with-current-input-port port . body)`  $\longrightarrow$  *value(s) of body* syntax  
`(with-current-output-port port . body)`  $\longrightarrow$  *value(s) of body* syntax  
`(with-error-output-port port . body)`  $\longrightarrow$  *value(s) of body* syntax

These special forms are simply syntactic sugar for the `with-current-input-port*` procedure and friends.



<code>(set-current-input-port! port)</code>	$\longrightarrow$ <i>undefined</i>	procedure
<code>(set-current-output-port! port)</code>	$\longrightarrow$ <i>undefined</i>	procedure
<code>(set-error-output-port! port)</code>	$\longrightarrow$ <i>undefined</i>	procedure

These procedures alter the dynamic binding of the current I/O port procedures to new values.

<code>(close fd/port)</code>	$\longrightarrow$ <i>boolean</i>	procedure
------------------------------	----------------------------------	-----------

Close the port or file descriptor.

If *fd/port* is a file descriptor, and it has a port allocated to it, the port is shifted to a new file descriptor created with `(dup fd/port)` before closing *fd/port*. The port then has its revealed count set to zero. This reflects the design criteria that ports are not associated with file descriptors, but with open files.

To close a file descriptor, and any associated port it might have, you must instead say one of (as appropriate):

```
(close (fdes->inport fd))
(close (fdes->outport fd))
```

The procedure returns true if it closed an open port. If the port was already closed, it returns false; this is not an error.

<code>(stdports-&gt;stdio)</code>	$\longrightarrow$ <i>undefined</i>	procedure
<code>(stdio-&gt;stdports)</code>	$\longrightarrow$ <i>undefined</i>	procedure

These two procedures are used to synchronise Unix' standard I/O file descriptors and Scheme's current I/O ports.

`(stdports->stdio)` causes the standard I/O file descriptors (0, 1, and 2) to take their values from the current I/O ports. It is exactly equivalent to the series of redirections:<sup>1</sup>

```
(dup (current-input-port) 0)
(dup (current-output-port) 1)
(dup (error-output-port) 2)
```

`stdio->stdports` causes the bindings of the current I/O ports to be changed to ports constructed over the standard I/O file descriptors. It is exactly equivalent to the series of assignments

```
(set-current-input-port! (fdes->inport 0))
(set-current-output-port! (fdes->outport 1))
(set-error-output-port! (fdes->outport 2))
```

---

<sup>1</sup>Why not `move->fdes`? Because the current output port and error port might be the same port.

However, you are more likely to find the dynamic-extent variant, `with-stdio-ports*`, below, to be of use in general programming.

(with-stdio-ports\* *thunk*)   → *value(s) of thunk*                    procedure  
 (with-stdio-ports *body...*)   → *value(s) of body*                    syntax

`with-stdio-ports*` binds the standard ports (`current-input-port`), (`current-output-port`), and (`error-output-port`) to be ports on file descriptors 0, 1, 2, and then calls *thunk*. It is equivalent to:

```
(with-current-input-port (fdes->inport 0)
  (with-current-output-port (fdes->inport 1)
    (with-error-output-port (fdes->outport 2)
      (thunk))))
```

The `with-stdio-ports` special form is merely syntactic sugar.

### 3.2.3 String ports

Scheme 48 has string ports, which you can use. Scsh has not committed to the particular interface or names that Scheme 48 uses, so be warned that the interface described herein may be liable to change.

(make-string-input-port *string*)   → *port*                            procedure

Returns a port that reads characters from the supplied string.

(make-string-output-port)   → *port*                                procedure  
 (string-output-port-output *port*)   → *string*                    procedure

A string output port is a port that collects the characters given to it into a string. The accumulated string is retrieved by applying `string-output-port-output` to the port.

(call-with-string-output-port *procedure*)   → *string*            procedure

The *procedure* value is called on a port. When it returns, `call-with-string-output-port` returns a string containing the characters that were written to that port during the execution of *procedure*.

### 3.2.4 Revealed ports and file descriptors

The material in this section and the following one is not critical for most applications. You may safely skim or completely skip this section on a first reading.

Dealing with Unix file descriptors in a Scheme environment is difficult. In Unix, open files are part of the process environment, and are referenced by

small integers called *file descriptors*. Open file descriptors are the fundamental way I/O redirections are passed to subprocesses, since file descriptors are preserved across fork's and exec's.

Scheme, on the other hand, uses ports for specifying I/O sources. Ports are garbage-collected Scheme objects, not integers. Ports can be garbage collected; when a port is collected, it is also closed. Because file descriptors are just integers, it's impossible to garbage collect them—you wouldn't be able to close file descriptor 3 unless there were no 3's in the system, and you could further prove that your program would never again compute a 3. This is difficult at best.

If a Scheme program only used Scheme ports, and never actually used file descriptors, this would not be a problem. But Scheme code must descend to the file descriptor level in at least two circumstances:

- when interfacing to foreign code
- when interfacing to a subprocess.

This causes a problem. Suppose we have a Scheme port constructed on top of file descriptor 2. We intend to fork off a program that will inherit this file descriptor. If we drop references to the port, the garbage collector may prematurely close file 2 before we fork the subprocess. The interface described below is intended to fix this and other problems arising from the mismatch between ports and file descriptors.

The Scheme kernel maintains a port table that maps a file descriptor to the Scheme port allocated for it (or, #f if there is no port allocated for this file descriptor). This is used to ensure that there is at most one open port for each open file descriptor.

The port data structure for file ports has two fields besides the descriptor: *revealed* and *closed?*. When a file port is closed with (`close port`), the port's file descriptor is closed, its entry in the port table is cleared, and the port's *closed?* field is set to true.

When a file descriptor is closed with (`close fdes`), any associated port is shifted to a new file descriptor created with (`dup fdes`). The port has its revealed count reset to zero (and hence becomes eligible for closing on GC). See discussion below. To really put a stake through a descriptor's heart without waiting for associated ports to be GC'd, you must say one of

```
(close (fdes->inport fdes))
(close (fdes->output fdes))
```

The *revealed* field is an aid to garbage collection. It is an integer semaphore. If it is zero, the port's file descriptor can be closed when the port is collected.

Essentially, the *revealed* field reflects whether or not the port's file descriptor has escaped to the Scheme user. If the Scheme user doesn't know what file descriptor is associated with a given port, then he can't possibly retain an "integer handle" on the port after dropping pointers to the port itself, so the garbage collector is free to close the file.

Ports allocated with `open-output-file` and `open-input-file` are unrevealed ports—i.e., *revealed* is initialised to 0. No one knows the port's file descriptor, so the file descriptor can be closed when the port is collected.

The functions `fdes->output-port`, `fdes->input-port`, `port->fdes` are used to shift back and forth between file descriptors and ports. When `port->fdes` reveals a port's file descriptor, it increments the port's *revealed* field. When the user is through with the file descriptor, he can call `(release-port-handle port)`, which decrements the count. The function `(call/fdes fd/port proc)` automates this protocol. `call/fdes` uses `dynamic-wind` to enforce the protocol. If *proc* throws out of the `call/fdes` application, the unwind handler releases the descriptor handle; if the user subsequently tries to throw *back* into *proc*'s context, the wind handler raises an error. When the user maps a file descriptor to a port with `fdes->output` or `fdes->input`, the port has its *revealed* field incremented.

Not all file descriptors are created by requests to make ports. Some are inherited on process invocation via `exec(2)`, and are simply part of the global environment. Subprocesses may depend upon them, so if a port is later allocated for these file descriptors, it should be considered as a revealed port. For example, when the Scheme shell's process starts up, it opens ports on file descriptors 0, 1, and 2 for the initial values of `(current-input-port)`, `(current-output-port)`, and `(error-output-port)`. These ports are initialised with *revealed* set to 1, so that `stdin`, `stdout`, and `stderr` are not closed even if the user drops the port.

Unrevealed file ports have the nice property that they can be closed when all pointers to the port are dropped. This can happen during gc, or at an `exec()`—since all memory is dropped at an `exec()`. No one knows the file descriptor associated with the port, so the `exec'd` process certainly can't refer to it.

This facility preserves the transparent close-on-collect property for file ports that are used in straightforward ways, yet allows access to the underlying Unix substrate without interference from the garbage collector. This is critical, since shell programming absolutely requires access to the Unix file descriptors, as their numerical values are a critical part of the process interface.

A port's underlying file descriptor can be shifted around with `dup(2)` when convenient. That is, the actual file descriptor on top of which a port is constructed can be shifted around underneath the port by the `scsh` kernel when

necessary. This is important, because when the user is setting up file descriptors prior to a `exec(2)`, he may explicitly use a file descriptor that has already been allocated to some port. In this case, the `scsh` kernel just shifts the port's file descriptor to some new location with `dup`, freeing up its old descriptor. This prevents errors from happening in the following scenario. Suppose we have a file open on port `f`. Now we want to run a program that reads input on file 0, writes output to file 1, errors to file 2, and logs execution information on file 3. We want to run this program with input from `f`. So we write:

```
(run (/usr/shivers/bin/prog)
  (> 1 output.txt)
  (> 2 error.log)
  (> 3 trace.log)
  (= 0 ,f))
```

Now, suppose by ill chance that, unbeknownst to us, when the operating system opened `f`'s file, it allocated descriptor 3 for it. If we blindly redirect `trace.log` into file descriptor 3, we'll clobber `f`! However, the port-shuffling machinery saves us: when the `run` form tries to `dup` `trace.log`'s file descriptor to 3, `dup` will notice that file descriptor 3 is already associated with an unrevealed port (*i.e.*, `f`). So, it will first move `f` to some other file descriptor. This keeps `f` alive and well so that it can subsequently be `dup`'d into descriptor 0 for `prog`'s `stdin`.

The port-shifting machinery makes the following guarantee: a port is only moved when the underlying file descriptor is closed, either by a `close()` or a `dup2()` operation. Otherwise a port/file-descriptor association is stable.

Under normal circumstances, all this machinery just works behind the scenes to keep things straightened out. The only time the user has to think about it is when he starts accessing file descriptors from ports, which he should almost never have to do. If a user starts asking what file descriptors have been allocated to what ports, he has to take responsibility for managing this information.

### 3.2.5 Port-mapping machinery

The procedures provided in this section are almost never needed. You may safely skim or completely skip this section on a first reading.

Here are the routines for manipulating ports in `scsh`. The important points to remember are:

- A file port is associated with an open file, not a particular file descriptor.

- The association between a file port and a particular file descriptor is never changed *except* when the file descriptor is explicitly closed. “Closing” includes being used as the target of a dup2, so the set of procedures below that close their targets are close, two-argument dup, and move->fdes. If the target file descriptor of one of these routines has an allocated port, the port will be shifted to another freshly-allocated file descriptor, and marked as unrevealed, thus preserving the port but freeing its old file descriptor.

These rules are what is necessary to “make things work out” with no surprises in the general case.

(fdes->inport <i>fd</i> )	→	<i>port</i>	procedure
(fdes->outport <i>fd</i> )	→	<i>port</i>	procedure
(port->fdes <i>port</i> )	→	<i>fixnum</i>	procedure

These increment the port’s revealed count.

(port-revealed <i>port</i> )	→	<i>integer or #f</i>	procedure
------------------------------	---	----------------------	-----------

Return the port’s revealed count if positive, otherwise #f.

(release-port-handle <i>port</i> )	→	<i>undefined</i>	procedure
------------------------------------	---	------------------	-----------

Decrement the port’s revealed count.

(call/fdes <i>fd/port consumer</i> )	→	<i>value(s) of consumer</i>	procedure
--------------------------------------	---	-----------------------------	-----------

Calls *consumer* on a file descriptor; takes care of revealed bookkeeping. If *fd/port* is a file descriptor, this is just (*consumer fd/port*). If *fd/port* is a port, calls *consumer* on its underlying file descriptor. While *consumer* is running, the port’s revealed count is incremented.

When call/fdes is called with port argument, you are not allowed to throw into *consumer* with a stored continuation, as that would violate the revealed-count bookkeeping.

(move->fdes <i>fd/port target-fd</i> )	→	<i>port or fdes</i>	procedure
--	---	---------------------	-----------

Maps fd→fd and port→port.

If *fd/port* is a file-descriptor not equal to *target-fd*, dup it to *target-fd* and close it. Returns *target-fd*.

If *fd/port* is a port, it is shifted to *target-fd*, by duping its underlying file-descriptor if necessary. *Fd/port*’s original file descriptor is closed (if it was different from *target-fd*). Returns the port. This operation resets *fd/port*’s revealed count to 1.

In all cases when *fd/port* is actually shifted, if there is a port already using *target-fd*, it is first relocated to some other file descriptor.

### 3.2.6 Unix I/O

<code>(dup <i>fd/port</i> [<i>newfd</i>])</code>	$\longrightarrow$ <i>fd/port</i>	procedure
<code>(dup-&gt;inport <i>fd/port</i> [<i>newfd</i>])</code>	$\longrightarrow$ <i>port</i>	procedure
<code>(dup-&gt;outport <i>fd/port</i> [<i>newfd</i>])</code>	$\longrightarrow$ <i>port</i>	procedure
<code>(dup-&gt;fdes <i>fd/port</i> [<i>newfd</i>])</code>	$\longrightarrow$ <i>fd</i>	procedure

These procedures provide the functionality of C's `dup()` and `dup2()`. The different routines return different types of values: `dup->inport`, `dup->outport`, and `dup->fdes` return input ports, output ports, and integer file descriptors, respectively. `dup`'s return value depends on the type of *fd/port*—it maps *fd*→*fd* and *port*→*port*.

These procedures use the Unix `dup()` syscall to replicate the file descriptor or file port *fd/port*. If a *newfd* file descriptor is given, it is used as the target of the `dup` operation, *i.e.*, the operation is a `dup2()`. In this case, procedures that return a port (such as `dup->inport`) will return one with the revealed count set to one. For example, `(dup (current-input-port) 5)` produces a new port with underlying file descriptor 5, whose revealed count is 1. If *newfd* is not specified, then the operating system chooses the file descriptor, and any returned port is marked as unrevealed.

If the *newfd* target is given, and some port is already using that file descriptor, the port is first quietly shifted (with another `dup`) to some other file descriptor (zeroing its revealed count).

Since Scheme doesn't provide read/write ports, `dup->inport` and `dup->outport` can be useful for getting an output version of an input port, or *vice versa*. For example, if *p* is an input port open on a tty, and we would like to do output to that tty, we can simply use `(dup->outport p)` to produce an equivalent output port for the tty. Be sure to open the file with the `open/read+write` flag for this.

<code>(seek <i>fd/port</i> <i>offset</i> [<i>whence</i>])</code>	$\longrightarrow$ <i>integer</i>	procedure
--	----------------------------------	-----------

Reposition the I/O cursor for a file descriptor or port. *whence* is one of `{seek/set, seek/delta, seek/end}`, and defaults to `seek/set`. If `seek/set`, then *offset* is an absolute index into the file; if `seek/delta`, then *offset* is a relative offset from the current I/O cursor; if `seek/end`, then *offset* is a relative offset from the end of file. The *fd/port* argument may be a port or an integer file descriptor. Not all such values are seekable; this is dependent on the OS implementation. The return value is the resulting position of the I/O cursor in the I/O stream.

*Oops:* The current implementation doesn't handle *offset* arguments that are not immediate integers (*i.e.*, representable in 30 bits).

*Oops:* The current implementation doesn't handle buffered ports.

(tell *fd/port*)  $\longrightarrow$  *integer* procedure

Returns the position of the I/O cursor in the the I/O stream. Not all file descriptors or ports support cursor-reporting; this is dependent on the OS implementation.

(open-file *fname flags [perms]*)  $\longrightarrow$  *port* procedure

*Perms* defaults to #o666. *Flags* is an integer bitmask, composed by or'ing together constants listed in table 3.1 (page 33). You must use exactly one of the open/read, open/write, or open/read+write flags. The returned port is an input port if the *flags* permit it, otherwise an output port. R5RS/Scheme 48/scsh do not have input/output ports, so it's one or the other. This should be fixed. (You can hack simultaneous I/O on a file by opening it r/w, taking the result input port, and duping it to an output port with dup->outputport.)

(open-input-file *fname [flags]*)  $\longrightarrow$  *port* procedure

(open-output-file *fname [flags perms]*)  $\longrightarrow$  *port* procedure

These are equivalent to open-file, after first setting the read/write bits of the *flags* argument to open/read or open/write, respectively. *Flags* defaults to zero for open-input-file, and

(bitwise-ior open/create open/truncate)

for open-output-file. These defaults make the procedures backwards-compatible with their unary R5RS definitions.

(open-fdes *fname flags [perms]*)  $\longrightarrow$  *integer* procedure

Returns a file descriptor.

(fdes-flags *fd/port*)  $\longrightarrow$  *integer* procedure

(set-fdes-flags *fd/port integer*)  $\longrightarrow$  *undefined* procedure

These procedures allow reading and writing of an open file's flags. The only such flag defined by POSIX is fdflags/close-on-exec; your Unix implementation may provide others.

These procedures should not be particularly useful to the programmer, as the scsh runtime already provides automatic control of the close-on-exec property. Unrevealed ports always have their file descriptors marked close-on-exec, as they can be closed when the scsh process execs a new program. Whenever the user reveals or unreveals a port's file descriptor, the runtime automatically sets or clears the flag for the programmer.



	Allowed operations	Status flag
<b>Open+Get+Set</b>	These flags can be used in open-file, fdes-status, and set-fdes-status calls.	open/append open/non-blocking open/async (Non-POSIX) open/fsync (Non-POSIX)
<b>Open+Get</b>	These flags can be used in open-file and fdes-status calls, but are ignored by set-fdes-status.	open/read open/write open/read+write open/access-mask
<b>Open</b>	These flags are only relevant in open-file calls; they are ignored by fdes-status and set-fdes-status calls.	open/create open/exclusive open/no-control-tty open/truncate

Table 3.1: Status flags for open-file, fdes-status and set-fdes-status. Only POSIX flags are guaranteed to be present; your operating system may define others. The open/access-mask value is not an actual flag, but a bit mask used to select the field for the open/read, open/write and open/read+write bits.

Programmers that manipulate this flag should be aware of these extra, automatic operations.

(fdes-status *fd/port*)  $\longrightarrow$  *integer* procedure  
(set-fdes-status *fd/port integer*)  $\longrightarrow$  *undefined* procedure

These procedures allow reading and writing of an open file’s status flags (table 3.1).

Note that this file-descriptor state is shared between file descriptors created by dup—if you create port *b* by applying dup to port *a*, and change *b*’s status flags, you will also have changed *a*’s status flags.

(pipe)  $\longrightarrow$  [*rport wport*] procedure

Returns two ports, the read and write end-points of a Unix pipe.

(read-string *nbytes [fd/port]*)  $\longrightarrow$  *string or #f* procedure  
(read-string! *str [fd/port start end]*)  $\longrightarrow$  *nread or #f* procedure

These calls read exactly as much data as you requested, unless there is not enough data (eof). `read-string!` reads the data into string *str* at the indices in the half-open interval  $[start, end)$ ; the default interval is the whole string:  $start = 0$  and  $end = (\text{string-length } string)$ . They will persistently retry on partial reads and when interrupted until (1) error, (2) eof, or (3) the input request is completely satisfied. Partial reads can occur when reading from an intermittent source, such as a pipe or tty.

`read-string` returns the string read; `read-string!` returns the number of characters read. They both return false at eof. A request to read zero bytes returns immediately, with no eof check.

The values of *start* and *end* must specify a well-defined interval in *str*, i.e.,  $0 \leq start \leq end \leq (\text{string-length } str)$ .

Any partially-read data is included in the error exception packet. Error returns on non-blocking input are considered an error.

`(read-string/partial nbytes [fd/port])`  $\longrightarrow$  *string or #f* procedure  
`(read-string!/partial str [fd/port start end])`  $\longrightarrow$  *nread or #f* procedure

These are atomic best-effort/forward-progress calls. Best effort: they may read less than you request if there is a lesser amount of data immediately available (e.g., because you are reading from a pipe or a tty). Forward progress: if no data is immediately available (e.g., empty pipe), they will block. Therefore, if you request an  $n > 0$  byte read, while you may not get everything you asked for, you will always get something (barring eof).

There is one case in which the forward-progress guarantee is cancelled: when the programmer explicitly sets the port to non-blocking I/O. In this case, if no data is immediately available, the procedure will not block, but will immediately return a zero-byte read.

`read-string/partial` reads the data into a freshly allocated string, which it returns as its value. `read-string!/partial` reads the data into string *str* at the indices in the half-open interval  $[start, end)$ ; the default interval is the whole string:  $start = 0$  and  $end = (\text{string-length } string)$ . The values of *start* and *end* must specify a well-defined interval in *str*, i.e.,  $0 \leq start \leq end \leq (\text{string-length } str)$ . It returns the number of bytes read.

A request to read zero bytes returns immediately, with no eof check.

In sum, there are only three ways you can get a zero-byte read: (1) you request one, (2) you turn on non-blocking I/O, or (3) you try to read at eof.

These are the routines to use for non-blocking input. They are also useful when you wish to efficiently process data in large blocks, and your algorithm is insensitive to the block size of any particular read operation.

<code>(select rvec wvec evec [timeout])</code>	$\longrightarrow$	<code>[rvec' wvec' evec']</code>	procedure
<code>(select! rvec wvec evec [timeout])</code>	$\longrightarrow$	<code>[nr nw ne]</code>	procedure

The `select` procedure allows a process to block and wait for events on multiple I/O channels. The `rvec` and `ewec` arguments are vectors of input ports and integer file descriptors; `wvec` is a vector of output ports and integer file descriptors. The procedure returns three vectors whose elements are subsets of the corresponding arguments. Every element of `rvec'` is ready for input; every element of `wvec'` is ready for output; every element of `ewec'` has an exceptional condition pending.

The `select` call will block until at least one of the I/O channels passed to it is ready for operation. For an input port this means that it either has data sitting its buffer or that the underlying file descriptor has data waiting. For an output port this means that it either has space available in the associated buffer or that the underlying file descriptor can accept output. For file descriptors, no buffers are checked, even if they have associated ports.

The `timeout` value can be used to force the call to time-out after a given number of seconds. It defaults to the special value `#f`, meaning wait indefinitely. A zero value can be used to poll the I/O channels.

If an I/O channel appears more than once in a given vector—perhaps occurring once as a Scheme port, and once as the port's underlying integer file descriptor—only one of these two references may appear in the returned vector. Buffered I/O ports are handled specially—if an input port's buffer is not empty, or an output port's buffer is not yet full, then these ports are immediately considered eligible for I/O without using the actual, primitive `select` system call to check the underlying file descriptor. This works pretty well for buffered input ports, but is a little problematic for buffered output ports.

The `select!` procedure is similar, but indicates the subset of active I/O channels by side-effecting the argument vectors. Non-active I/O channels in the argument vectors are overwritten with `#f` values. The call returns the number of active elements remaining in each vector. As a convenience, the vectors passed in to `select!` are allowed to contain `#f` values as well as integers and ports.

*Remark:* `Select` and `select!` do not call their POSIX counterparts directly—there is a POSIX `select` sitting at the very heart of the Scheme 48/scsh I/O system, so *all* multiplexed I/O is really `select`-based. Therefore, you cannot expect a performance increase from writing a single-threaded program using `select` and `select!` instead of writing a multi-threaded program where each thread handles one I/O connection.

The moral of this story is that `select` and `select!` make sense in only two situations: legacy code written for an older version of `scsh`, and programs which make inherent use of `select/select!` which do not benefit from multiple threads. Examples are network clients that send requests to multiple alternate servers and discard all but one of them. In any case, the `select-ports` and `select-port-channels` procedures described below are usually a preferable alternative to `select/select!`: they are much simpler to use, and also have a slightly more efficient implementation.

`(select-ports timeout port ...)`  $\longrightarrow$  *ready-ports* procedure

The `select-ports` call will block until at least one of the ports passed to it is ready for operation or until the timeout has expired. For an input port this means that it either has data sitting its buffer or that the underlying file descriptor has data waiting. For an output port this means that it either has space available in the associated buffer or that the underlying file descriptor can accept output.

The *timeout* value can be used to force the call to time out after a given number of seconds. A value of `#f` means to wait indefinitely. A zero value can be used to poll the ports.

`Select-ports` returns a list of the ports ready for operation. Note that this list may be empty if the timeout expired before any ports became ready.

`(select-port-channels timeout port ...)`  $\longrightarrow$  *ready-ports* procedure

`Select-port-channels` is like `select-ports`, except that it only looks at the operating system objects the ports refer to, ignoring any buffering performed by the ports.

*Remark:* `Select-port-channels` should be used with care: for example, if an input port has data in the buffer but no data available on the underlying file descriptor, `select-port-channels` will block, even though a read operation on the port would be able to complete without blocking.

`Select-port-channels` is intended for situations where the program is not checking for available data, but rather for waiting until a port has established a connection—for example, to a network port.

`(write-string string [fd/port start end])`  $\longrightarrow$  *undefined* procedure

This procedure writes all the data requested. If the procedure cannot perform the write with a single kernel call (due to interrupts or partial writes), it will perform multiple write operations until all the data is written or an error has occurred. A non-blocking I/O error is considered an

error. (Error exception packets for this syscall include the amount of data partially transferred before the error occurred.)

The data written are the characters of *string* in the half-open interval  $[start, end)$ . The default interval is the whole string:  $start = 0$  and  $end = (\text{string-length } string)$ . The values of *start* and *end* must specify a well-defined interval in *str*, i.e.,  $0 \leq start \leq end \leq (\text{string-length } str)$ . A zero-byte write returns immediately, with no error.

Output to buffered ports: *write-string*'s efforts end as soon as all the data has been placed in the output buffer. Errors and true output may not happen until a later time, of course.

`(write-string/partial string [fd/port start end])`  $\longrightarrow$  *nwritten* procedure

This routine is the atomic best-effort/forward-progress analog to *write-string*. It returns the number of bytes written, which may be less than you asked for. Partial writes can occur when (1) we write off the physical end of the media, (2) the write is interrupted, or (3) the file descriptor is set for non-blocking I/O.

If the file descriptor is not set up for non-blocking I/O, then a successful return from these procedures makes a forward progress guarantee—that is, a partial write took place of at least one byte:

- If we are at the end of physical media, and no write takes place, an error exception is raised. So a return implies we wrote *something*.
- If the call is interrupted after a partial transfer, it returns immediately. But if the call is interrupted before any data transfer, then the write is retried.

If we request a zero-byte write, then the call immediately returns 0. If the file descriptor is set for non-blocking I/O, then the call may return 0 if it was unable to immediately write anything (e.g., full pipe). Barring these two cases, a write either returns *nwritten*  $> 0$ , or raises an error exception.

Non-blocking I/O is only available on file descriptors and unbuffered ports. Doing non-blocking I/O to a buffered port is not well-defined, and is an error (the problem is the subsequent flush operation).

*Oops:* *write-string/partial* is currently not implemented. Consider using threads to achieve the same functionality.

### 3.2.7 Buffered I/O

Scheme 48 ports use buffered I/O—data is transferred to or from the OS in blocks. Scsh provides control of this mechanism: the programmer may force saved-up output data to be transferred to the OS when he chooses, and may also choose which I/O buffering policy to employ for a given port (or turn buffering off completely).

It can be useful to turn I/O buffering off in some cases, for example when an I/O stream is to be shared by multiple subprocesses. For this reason, scsh allocates an unbuffered port for file descriptor 0 at start-up time. Because shells frequently share stdin with subprocesses, if the shell does buffered reads, it might “steal” input intended for a subprocess. For this reason, all shells, including sh, csh, and scsh, read stdin unbuffered. Applications that can tolerate buffered input on stdin can reset (`current-input-port`) to block buffering for higher performance.

{Note So support `peek-char` a Scheme implementation has to maintain a buffer for all input ports. In scsh, for “unbuffered” input ports the buffer size is one. As you cannot request less than one character there is no unrequested reading so this can still be called “unbuffered input”.}

(`set-port-buffering port policy [size]`)  $\longrightarrow$  *undefined* procedure

This procedure allows the programmer to assign a particular I/O buffering policy to a port, and to choose the size of the associated buffer. It may only be used on new ports, *i.e.*, before I/O is performed on the port. There are three buffering policies that may be chosen:

<code>bufpol/block</code>	General block buffering (general default)
<code>bufpol/line</code>	Line buffering (tty default)
<code>bufpol/none</code>	Direct I/O—no buffering <sup>2</sup>

The line buffering policy flushes output whenever a newline is output; whenever the buffer is full; or whenever an input is read from stdin. Line buffering is the default for ports open on terminal devices.

*Oops:* The current implementation doesn’t support `bufpol/line`.

The *size* argument requests an I/O buffer of *size* bytes. For output ports, *size* must be non-negative, for input ports *size* must be positive. If not given, a reasonable default is used. For output ports, if given and zero, buffering is turned off (*i.e.*, *size* = 0 for any policy is equivalent to *policy* = `bufpol/none`). For input ports, setting the size to one corresponds to unbuffered input as defined above. If given, *size* must be zero respectively one for `bufpol/none`.

(`force-output [fd/port]`)  $\longrightarrow$  *undefined* procedure

This procedure does nothing when applied to an integer file descriptor or unbuffered port. It flushes buffered output when applied to a buffered port, and raises a write-error exception on error. Returns no value.

`(flush-all-ports)`  $\longrightarrow$  *undefined* procedure

This procedure flushes all open output ports with buffered data.

### 3.2.8 File locking

Scsh provides POSIX advisory file locking. *Advisory* locks are locks that can be checked by user code, but do not affect other I/O operations. For example, if a process has an exclusive lock on a region of a file, other processes will not be able to obtain locks on that region of the file, but they will still be able to read and write the file with no hindrance. Using advisory locks requires cooperation amongst the agents accessing the shared resource.

*Remark:* Unfortunately, POSIX file locks are associated with actual files, not with associated open file descriptors. Once a process locks a file, using some file descriptor *fd*, the next time *any* file descriptor referencing that file is closed, all associated locks are released. This severely limits the utility of POSIX advisory file locks, and we'd recommend caution when using them. It is not without reason that the FreeBSD man pages refer to POSIX file locking as "completely stupid."

Scsh moves Scheme ports from file descriptor to file descriptor with `dup()` and `close()` as required by the runtime, so it is impossible to keep file locks open across one of these shifts. Hence we can only offer POSIX advisory file locking directly on raw integer file descriptors; regrettably, there are no facilities for locking Scheme ports.

Note that once a Scheme port is revealed in scsh, the runtime will not shift the port around with `dup()` and `close()`. This means the file-locking procedures can then be applied to the port's associated file descriptor.

POSIX allows the user to lock a region of a file with either an exclusive or shared lock. Locked regions are described by the *lock-region* record:

```
(define-record lock-region
  exclusive?
  start
  len
  whence
  proc)
```

- The `exclusive?` field is true if the lock is exclusive; false if it is shared.
- The `whence` field is one of the values from the `seek` call: `seek/set`, `seek/delta`, or `seek/end`, and determines the interpretation of the `start` field:
  - If `seek/set`, the `start` value is simply an absolute index into the file.
  - If `seek/delta`, the `start` value is an offset from the file descriptor's current position in the file.
  - If `seek/end`, the `start` value is an offset from the end of the file.

The region of the file being locked is given by the `start` and `len` fields; if `len` is zero, it means "infinity," that is, the region extends from the starting point through the end of the file, even as the file is extended by subsequent write operations.

- The `proc` field gives the process object for the process holding the region lock, when relevant (see `get-lock-region` below).

`(make-lock-region exclusive? start len [whence])`  $\longrightarrow$  *lock-region* procedure

This procedure makes a lock-region record. The `whence` field defaults to `seek/set`.

`(lock-region fdes lock)`  $\longrightarrow$  *undefined* procedure

`(lock-region/no-block fdes lock)`  $\longrightarrow$  *boolean* procedure

These procedures lock a region of the file referenced by file descriptor *fdes*. The `lock-region` procedure blocks until the lock is granted; the non-blocking variant returns a boolean indicating whether or not the lock was granted. To take an exclusive (write) lock, you must have the file descriptor open with write access; to take a shared (read) lock, you must have the file descriptor open with read access.

`(get-lock-region fdes lock)`  $\longrightarrow$  *lock-region or #f* procedure

Return the first lock region on *fdes* that would conflict with lock region *lock*. If there is no such lock region, return `false`. This procedure fills out the `proc` field of the returned lock region, and is the only procedure that has anything to do with this field. (See section 3.4.1 for a description of process objects.) Note that if you apply this procedure to a file system that is shared across multiple operating systems (*i.e.*, an NFS file system), the `proc` field may be ambiguous. We note, again, that POSIX advisory file locking is not a terribly useful or well-designed facility.



(unlock-region *fdes lock*)  $\longrightarrow$  *undefined* procedure  
 Release a lock from a file.

(with-region-lock\* *fdes lock thunk*)  $\longrightarrow$  *value(s) of thunk* procedure  
 (with-region-lock *fdes lock body ...*)  $\longrightarrow$  *value(s) of body* syntax

This procedure obtains the requested lock, and then calls (*thunk*). When *thunk* returns, the lock is released. A non-local exit (*e.g.*, throwing to a saved continuation or raising an exception) also causes the lock to be released.

After a normal return from *thunk*, its return values are returned by *with-region-lock\**. The *with-region-lock* special form is equivalent syntactic sugar.

### 3.3 File system

Besides the following procedures, which allow access to the computer's file system, *scsh* also provides a set of procedures which manipulate file *names*. These string-processing procedures are documented in section 5.1.

(create-directory *fname [perms override?]*)  $\longrightarrow$  *undefined* procedure  
 (create-fifo *fname [perms override?]*)  $\longrightarrow$  *undefined* procedure  
 (create-hard-link *oldname newname [override?]*)  $\longrightarrow$  *undefined* procedure  
 (create-symlink *old-name new-name [override?]*)  $\longrightarrow$  *undefined* procedure

These procedures create objects of various kinds in the file system.

The *override?* argument controls the action if there is already an object in the file system with the new name:

#f        signal an error (default)  
 'query   prompt the user  
*other*    delete the old object (with *delete-file* or *delete-directory*, as appropriate) before creating the new object.

*Perms* defaults to #o777 (but is masked by the current *umask*).

*Remark:* Currently, if you try to create a hard or symbolic link from a file to itself, you will error out with *override?* false, and simply delete your file with *override?* true. Catching this will require some sort of *true-name* procedure, which I currently do not have.

```

(delete-directory fname) → undefined           procedure
(delete-file fname) → undefined               procedure
(delete-filesys-object fname) → undefined       procedure

```

These procedures delete objects from the file system. The `delete-filesys-object` procedure will delete an object of any type from the file system: files, (empty) directories, symlinks, fifos, *etc.*.

If the object being deleted doesn't exist, `delete-directory` and `delete-file` raise an error, while `delete-filesys-object` simply returns.

```

(read-symlink fname) → string                 procedure

```

Return the filename referenced by symbolic link *fname*.

```

(rename-file old-fname new-fname [override?]) → undefined procedure

```

If you override an existing object, then *old-fname* and *new-fname* must type-match—either both directories, or both non-directories. This is required by the semantics of Unix `rename()`.

*Remark:* There is an unfortunate atomicity problem with the `rename-file` procedure: if you specify `no-override`, but create file *new-fname* sometime between `rename-file`'s existence check and the actual rename operation, your file will be clobbered with *old-fname*. There is no way to fix this problem, given the semantics of Unix `rename()`; at least it is highly unlikely to occur in practice.

```

(set-file-mode fname/fd/port mode) → undefined   procedure
(set-file-owner fname/fd/port uid) → undefined   procedure
(set-file-group fname/fd/port gid) → undefined   procedure

```

These procedures set the permission bits, owner id, and group id of a file, respectively. The file can be specified by giving the file name, or either an integer file descriptor or a port open on the file. Setting file user ownership usually requires root privileges.

```

(set-file-times fname [access-time mod-time]) → undefined procedure

```

This procedure sets the access and modified times for the file *fname* to the supplied values (see section 3.10 for the `scsh` representation of time). If neither time argument is supplied, they are both taken to be the current time. You must provide both times or neither. If the procedure completes successfully, the file's time of last status-change (`ctime`) is set to the current time.

(sync-file *fd/port*) → *undefined* procedure  
 (sync-file-system) → *undefined* procedure

Calling `sync-file` causes Unix to update the disk data structures for a given file. If *fd/port* is a port, any buffered data it may have is first flushed. Calling `sync-file-system` synchronises the kernel's entire file system with the disk.

These procedures are not POSIX. Interestingly enough, `sync-file-system` doesn't actually do what it is claimed to do. We just threw it in for humor value. See the `sync(2)` man page for Unix enlightenment.

(truncate-file *fname/fd/port len*) → *undefined* procedure

The specified file is truncated to *len* bytes in length.

(file-info *fname/fd/port [chase?]*) → *file-info-record* procedure

The `file-info` procedure returns a record structure containing everything there is to know about a file. If the *chase?* flag is true (the default), then the procedure chases symlinks and reports on the files to which they refer. If *chase?* is false, then the procedure checks the actual file itself, even if it's a symlink. The *chase?* flag is ignored if the file argument is a file descriptor or port.

The value returned is a *file-info record*, defined to have the following structure:

```
(define-record file-info
  type      ; {block-special, char-special, directory,
             ;      fifo, regular, socket, symlink}
  device    ; Device file resides on.
  inode     ; File's inode.
  mode      ; File's mode bits: permissions, setuid, setgid
  nlinks    ; Number of hard links to this file.
  uid       ; Owner of file.
  gid       ; File's group id.
  size      ; Size of file, in bytes.
  atime     ; Time of last access.
  mtime     ; Time of last mod.
  ctime)    ; Time of last status change.
```

The `uid` field of a `file-info` record is accessed with the procedure

```
(file-info:uid x)
```

and similarly for the other fields. The type field is a symbol; all other fields are integers. A file-info record is discriminated with the `file-info?` predicate.

The following procedures all return selected information about a file; they are built on top of `file-info`, and are called with the same arguments that are passed to it.

Procedure	returns
<code>file-type</code>	type
<code>file-inode</code>	inode
<code>file-mode</code>	mode
<code>file-nlinks</code>	nlinks
<code>file-owner</code>	uid
<code>file-group</code>	gid
<code>file-size</code>	size
<code>file-last-access</code>	atime
<code>file-last-mod</code>	mtime
<code>file-last-status-change</code>	ctime

Example:

```
;; All my files in /usr/tmp:
(filter (lambda (f) (= (file-owner f) (user-uid)))
        (directory-files "/usr/tmp"))
```

*Remark:* `file-info` was named `file-attributes` in releases of `scsh` prior to release 0.4. We changed the name to `file-info` for consistency with the other information-retrieval procedures in `scsh`: `user-info`, `group-info`, `host-info`, `network-info`, `service-info`, and `protocol-info`.

The `file-attributes` binding is still supported in the current release of `scsh`, but is deprecated, and may go away in a future release.

<code>(file-directory? fname/fd/port [chase?])</code>	$\rightarrow$ <i>boolean</i>	procedure
<code>(file-fifo? fname/fd/port [chase?])</code>	$\rightarrow$ <i>boolean</i>	procedure
<code>(file-regular? fname/fd/port [chase?])</code>	$\rightarrow$ <i>boolean</i>	procedure
<code>(file-socket? fname/fd/port [chase?])</code>	$\rightarrow$ <i>boolean</i>	procedure
<code>(file-special? fname/fd/port [chase?])</code>	$\rightarrow$ <i>boolean</i>	procedure
<code>(file-symlink? fname/fd/port)</code>	$\rightarrow$ <i>boolean</i>	procedure

These procedures are file-type predicates that test the type of a given file. They are applied to the same arguments to which `file-info` is applied; the sole exception is `file-symlink?`, which does not take the optional *chase?* second argument.

For example,

```
(file-directory? "/usr/dalbertz")  ==>  #t
```

There are variants of these procedures which work directly on file-info records:

(file-info-directory? <i>file-info</i> )	→ <i>boolean</i>	procedure
(file-info-fifo? <i>file-info</i> )	→ <i>boolean</i>	procedure
(file-info-regular? <i>file-info</i> )	→ <i>boolean</i>	procedure
(file-info-socket? <i>file-info</i> )	→ <i>boolean</i>	procedure
(file-info-special? <i>file-info</i> )	→ <i>boolean</i>	procedure
(file-info-symlink? <i>file-info</i> )	→ <i>boolean</i>	procedure

The following set of procedures are a convenient means to work on the permission bits of a file:

(file-not-readable? <i>fname/fd/port</i> )	→ <i>boolean</i>	procedure
(file-not-writable? <i>fname/fd/port</i> )	→ <i>boolean</i>	procedure
(file-not-executable? <i>fname/fd/port</i> )	→ <i>boolean</i>	procedure

Returns:

Value	meaning
#f	Access permitted
'search-denied	Can't stat—a protected directory is blocking access.
'permission	Permission denied.
'no-directory	Some directory doesn't exist.
'nonexistent	File doesn't exist.

A file is considered writeable if either (1) it exists and is writeable or (2) it doesn't exist and the directory is writeable. Since symlink permission bits are ignored by the filesystem, these calls do not take a *chase?* flag.

Note that these procedures use the process' *effective* user and group ids for permission checking. POSIX defines an `access()` function that uses the process' real uid and gids. This is handy for setuid programs that would like to find out if the actual user has specific rights; scsh ought to provide this functionality (but doesn't at the current time).

There are several problems with these procedures. First, there's an atom-icity issue. In between checking permissions for a file and then trying an operation on the file, another process could change the permissions, so a return value from these functions guarantees nothing. Second, the code special-cases permission checking when the uid is root—if the file exists,

root is assumed to have the requested permission. However, not even root can write a file that is on a read-only file system, such as a CD ROM. In this case, `file-not-writable?` will lie, saying that root has write access, when in fact the opening the file for write access will fail. Finally, write permission confounds write access and create access. These should be disentangled.

Some of these problems could be avoided if POSIX had a real-uid variant of the `access()` call we could use, but the atomicity issue is still a problem. In the final analysis, the only way to find out if you have the right to perform an operation on a file is to try and open it for the desired operation. These permission-checking functions are mostly intended for script-writing, where loose guarantees are tolerated.

<code>(file-readable? fname/fd/port)</code>	$\longrightarrow$ <i>boolean</i>	procedure
<code>(file-writable? fname/fd/port)</code>	$\longrightarrow$ <i>boolean</i>	procedure
<code>(file-executable? fname/fd/port)</code>	$\longrightarrow$ <i>boolean</i>	procedure

These procedures are the logical negation of the preceding `file-not-...?` procedures. Refer to them for a discussion of their problems and limitations.

<code>(file-info-not-readable? file-info)</code>	$\longrightarrow$ <i>boolean</i>	procedure
<code>(file-info-not-writable? file-info)</code>	$\longrightarrow$ <i>boolean</i>	procedure
<code>(file-info-not-executable? file-info)</code>	$\longrightarrow$ <i>boolean</i>	procedure
<code>(file-info-readable? file-info)</code>	$\longrightarrow$ <i>boolean</i>	procedure
<code>(file-info-writable? file-info)</code>	$\longrightarrow$ <i>boolean</i>	procedure
<code>(file-info-executable? file-info)</code>	$\longrightarrow$ <i>boolean</i>	procedure

There are variants which work directly on `file-info` records.

<code>(file-not-exists? fname/fd/port [chase?])</code>	$\longrightarrow$ <i>object</i>	procedure
--	---------------------------------	-----------

Returns:

<code>#f</code>	Exists.
<code>#t</code>	Doesn't exist.
<code>'search-denied</code>	Some protected directory is blocking the search.

<code>(file-exists? fname/fd/port [chase?])</code>	$\longrightarrow$ <i>boolean</i>	procedure
--	----------------------------------	-----------

This is simply `(not (file-not-exists? fname [chase?]))`

<code>(directory-files [dir dotfiles?])</code>	$\longrightarrow$ <i>string list</i>	procedure
--	--------------------------------------	-----------

Return the list of files in directory *dir*, which defaults to the current working directory. The *dotfiles?* flag (default *#f*) causes dot files to be included in the list. Regardless of the value of *dotfiles?*, the two files *.* and *..* are *never* returned.

The directory *dir* is not prepended to each file name in the result list. That is,

```
(directory-files "/etc")
```

returns

```
("chown" "exports" "fstab" ...)
```

*not*

```
("etc/chown" "etc/exports" "etc/fstab" ...)
```

To use the files in returned list, the programmer can either manually prepend the directory:

```
(map (λ (f) (string-append dir "/" f)) files)
```

or *cd* to the directory before using the file names:

```
(with-cwd dir
  (for-each delete-file (directory-files)))
```

or use the *glob* procedure, defined below.

A directory list can be generated by `(run/strings (ls))`, but this is unreliable, as filenames with whitespace in their names will be split into separate entries. Using *directory-files* is reliable.

```
(open-directory-stream dir) → directory-stream-record    procedure
(read-directory-stream directory-stream-record) → string or #f procedure
(close-directory-stream directory-stream-record) → undefined procedure
```

These functions implement a direct interface to the `opendir()/readdir()/closedir()` family of functions for processing directory streams. `(open-directory-stream dir)` creates a stream of files in the directory *dir*. `(read-directory-stream directory-stream)` returns the next file in the stream or *#f* if no such file exists. Finally, `(close-directory-stream directory-stream)` closes the stream.

```
(glob pat1 ...) → string list    procedure
```

Glob each pattern against the filesystem and return the sorted list. Duplicates are not removed. Patterns matching nothing are not included literally.<sup>3</sup> C shell `{a,b,c}` patterns are expanded. Backslash quotes characters, turning off the special meaning of `{`, `}`, `*`, `[`, `]`, and `?`.

---

<sup>3</sup>Why bother to mention such a silly possibility? Because that is what *sh* does.

Note that the rules of backslash for Scheme strings and glob patterns work together to require four backslashes in a row to specify a single literal backslash. Fortunately, it is very rare that a backslash occurs in a Unix file name.

A glob subpattern will not match against dot files unless the first character of the subpattern is a literal `"."`. Further, a dot subpattern will not match the files `.` or `..` unless it is a constant pattern, as in `(glob "../*/*.c")`. So a directory's dot files can be reliably generated with the simple glob pattern `"*"`.

Some examples:

```
(glob "*.c" "*.h")
;; All the C and #include files in my directory.

(glob "*.c" "*/*.c")
;; All the C files in this directory and
;; its immediate subdirectories.

(glob "lexer/*.c" "parser/*.c")
(glob "{lexer,parser}/*.c")
;; All the C files in the lexer and parser dirs.

(glob "\\{lexer,parser}\\/*.c")
;; All the C files in the strange
;; directory "{lexer,parser}".

(glob "*\\*")
;; All the files ending in "*", e.g.
;; ("foo*" "bar*")

(glob "*lexer*")
("mylexer.c" "lexer1.notes")
;; All files containing the string "lexer".

(glob "lexer")
;; Either ("lexer") or ().
```

If the first character of the pattern (after expanding braces) is a slash, the search begins at root; otherwise, the search begins in the current working directory.

If the last character of the pattern (after expanding braces) is a slash, then the result matches must be directories, *e.g.*,



```
(glob "/usr/man/man?/") ==>
  ("/usr/man/man1/" "/usr/man/man2/" ...)
```

Globbering can sometimes be useful when we need a list of a directory's files where each element in the list includes the pathname for the file. Compare:

```
(directory-files "../include") ==>
  ("cig.h" "decls.h" ...)

(glob "../include/*") ==>
  ("../include/cig.h" "../include/decls.h" ...)
```

`(glob-quote str)`  $\longrightarrow$  *string* procedure

Returns a constant glob pattern that exactly matches *str*. All wild-card characters in *str* are quoted with a backslash.

```
(glob-quote "Any *.c files?")
==> "Any \*.c files\?"
```

`(file-match root dot-files? pat1 pat2 ... patn)`  $\longrightarrow$  *string list* procedure

{Note This procedure is deprecated, and will probably either go away or be substantially altered in a future release. New code should not call this procedure. The problem is that it relies upon Posix-notation regular expressions; the rest of scsh has been converted over to the new SRE notation.}

`file-match` provides a more powerful file-matching service, at the expense of a less convenient notation. It is intermediate in power between most shell matching machinery and recursive `find(1)`.

Each pattern is a regexp. The procedure searches from *root*, matching the first-level files against pattern *pat*<sub>1</sub>, the second-level files against *pat*<sub>2</sub>, and so forth. The list of files matching the whole path pattern is returned, in sorted order. The matcher uses Spencer's regular expression package.

The files `.` and `..` are never matched. Other dot files are only matched if the *dot-files?* argument is `#t`.

A given *pat*<sub>*i*</sub> pattern is matched as a regexp, so it is not forced to match the entire file name. *E.g.*, pattern `"t"` matches any file containing a `"t"` in its name, while pattern `"^t$"` matches only a file whose entire name is `"t"`.

The *pat*<sub>*i*</sub> patterns can be more general than stated above.

- A single pattern can specify multiple levels of the path by embedding / characters within the pattern. For example, the pattern "a/b/c" gives a match equivalent to the list of patterns "a" "b" "c".
- A  $pat_i$  pattern can be a procedure, which is used as a match predicate. It will be repeatedly called with a candidate file-name to test. The file-name will be the entire path accumulated. If the procedure raises an error condition, file-match will catch the error and treat it as a failed match. This keeps file-match from being blown out of the water by applying tests to dangling symlinks and other similar situations.

Some examples:

```
(file-match "/usr/lib" #f "m$" "^tab") =>
  ("/usr/lib/term/tab300" "/usr/lib/term/tab300-12" ...)
```

```
(file-match "." #f "^lex|parse|codegen$" "\\..c$") =>
  ("lex/lex.c" "lex/lexinit.c" "lex/test.c"
   "parse/actions.c" "parse/error.c" "parse/test.c"
   "codegen/io.c" "codegen/walk.c")
```

```
(file-match "." #f "^lex|parse|codegen$/\\..c$")
;; The same.
```

```
(file-match "." #f file-directory?)
;; Return all subdirs of the current directory.
```

```
(file-match "/" #f file-directory?) =>
  ("/bin" "/dev" "/etc" "/tmp" "/usr")
;; All subdirs of root.
```

```
(file-match "." #f "\\..c")
;; All the C files in my directory.
```

```
(define (ext extension)
  (λ (fn) (string-suffix? fn extension)))
```

```
(define (true . x) #t)
```

```

(file-match "." #f ".\\\.c")
(file-match "." #f "" "\\\.c")
(file-match "." #f true "\\\.c")
(file-match "." #f true (ext "c"))
;; All the C files of all my immediate subdirs.

(file-match "." #f "lexer") =>
  ("mylexer.c" "lexer.notes")
;; Compare with (glob "lexer"), above.

```

Note that when *root* is the current working directory ((".")), when it is converted to directory form, it becomes "", and doesn't show up in the result file-names.

It is regrettable that the regexp wild card char, ".", is such an important file name literal, as dot-file prefix and extension delimiter.

`(create-temp-file [prefix])`  $\longrightarrow$  *string* procedure

`Create-temp-file` creates a new temporary file and return its name. The optional argument specifies the filename prefix to use, and defaults to the value of "\$TMPDIR/*pid*" if \$TMPDIR is set and to "/var/tmp/*pid*" otherwise, where *pid* is the current process' id. The procedure generates a sequence of filenames that have *prefix* as a common prefix, looking for a filename that doesn't already exist in the file system. When it finds one, it creates it, with permission #o600 and returns the filename. (The file permission can be changed to a more permissive permission with `set-file-mode` after being created).

This file is guaranteed to be brand new. No other process will have it open. This procedure does not simply return a filename that is very likely to be unused. It returns a filename that definitely did not exist at the moment `create-temp-file` created it.

It is not necessary for the process' pid to be a part of the filename for the uniqueness guarantees to hold. The pid component of the default prefix simply serves to scatter the name searches into sparse regions, so that collisions are less likely to occur. This speeds things up, but does not affect correctness.

Security note: doing I/O to files created this way in /var/tmp/ is not necessarily secure. General users have write access to /var/tmp/, so even if an attacker cannot access the new temp file, he can delete it and replace it with one of his own. A subsequent open of this filename will then give you his file, to which he has access rights. There are several ways to defeat this attack,

1. Use `temp-file-iterate`, below, to return the file descriptor allocated when the file is opened. This will work if the file only needs to be opened once.
2. If the file needs to be opened twice or more, create it in a protected directory, *e.g.*, `$HOME`.
3. Ensure that `/var/tmp` has its sticky bit set. This requires system administrator privileges.

The actual default prefix used is controlled by the dynamic variable `*temp-file-template*`, and can be overridden for increased security. See `temp-file-iterate`.

```
(temp-file-iterate maker [template])  →  object+           procedure
*temp-file-template*                  string
```

This procedure can be used to perform certain atomic transactions on the file system involving filenames. Some examples:

- Linking a file to a fresh backup temp name.
- Creating and opening an unused, secure temp file.
- Creating an unused temporary directory.

This procedure uses *template* to generate a series of trial file names. *Template* is a format control string, and defaults to

```
"$TMPDIR/pid.~a"
```

if `$TMPDIR` is set and

```
"/var/tmp/pid.~a"
```

otherwise where *pid* is the current process' process id. File names are generated by calling `format` to instantiate the template's `~a` field with a varying string.

*Maker* is a procedure which is serially called on each file name generated. It must return at least one value; it may return multiple values. If the first return value is `#f` or if *maker* raises the `errno/exist` `errno` exception, `temp-file-iterate` will loop, generating a new file name and calling *maker* again. If the first return value is true, the loop is terminated, returning whatever value(s) *maker* returned.

After a number of unsuccessful trials, `temp-file-iterate` may give up and signal an error.

Thus, if we ignore its optional *prefix* argument, `create-temp-file` could be defined as:

```

(define (create-temp-file)
  (let ((flags (bitwise-ior open/create open/exclusive)))
    (temp-file-iterate
      (lambda (f)
        (close (open-output-file f flags #o600))
        f))))

```

To rename a file to a temporary name:

```

(temp-file-iterate (lambda (backup)
  (create-hard-link old-file backup)
  backup)
  ".#temp.~a") ; Keep link in cwd.
(delete-file old-file)

```

Recall that `scsh` reports syscall failure by raising an error exception, not by returning an error code. This is critical to this example—the programmer can assume that if the `temp-file-iterate` call returns, it returns successfully. So the following `delete-file` call can be reliably invoked, safe in the knowledge that the backup link has definitely been established.

To create a unique temporary directory:

```

(temp-file-iterate (lambda (dir) (create-directory dir) dir)
  "/var/tmp/tempdir.~a")

```

Similar operations can be used to generate unique symlinks and fifos, or to return values other than the new filename (*e.g.*, an open file descriptor or port).

The default template is in fact taken from the value of the dynamic variable `*temp-file-template*`, which itself defaults to `"$TMPDIR/pid.~a"` if `$TMPDIR` is set and `"/usr/tmp/pid.~a"` otherwise, where *pid* is the `scsh` process' pid. For increased security, a user may wish to change the template to use a directory not allowing world write access (*e.g.*, his home directory).

`(temp-file-channel)`  $\longrightarrow$  *[inp outp]* procedure

This procedure can be used to provide an interprocess communications channel with arbitrary-sized buffering. It returns two values, an input port and an output port, both open on a new temp file. The temp file itself is deleted from the Unix file tree before `temp-file-channel` returns, so the file is essentially unnamed, and its disk storage is reclaimed as soon as the two ports are closed.

`Temp-file-channel` is analogous to `port-pipe` with two exceptions:

- If the writer process gets ahead of the reader process, it will not hang waiting for some small pipe buffer to drain. It will simply buffer the data on disk. This is good.
- If the reader process gets ahead of the writer process, it will also not hang waiting for data from the writer process. It will simply see and report an end of file. This is bad.

In order to ensure that an end-of-file returned to the reader is legitimate, the reader and writer must serialise their I/O. The simplest way to do this is for the reader to delay doing input until the writer has completely finished doing output, or exited.

### 3.4 Processes

<code>(exec prog arg<sub>1</sub> ... arg<sub>n</sub>)</code>	$\longrightarrow$	<i>no return value</i>	procedure
<code>(exec-path prog arg<sub>1</sub> ... arg<sub>n</sub>)</code>	$\longrightarrow$	<i>no return value</i>	procedure
<code>(exec/env prog env arg<sub>1</sub> ... arg<sub>n</sub>)</code>	$\longrightarrow$	<i>no return value</i>	procedure
<code>(exec-path/env prog env arg<sub>1</sub> ... arg<sub>n</sub>)</code>	$\longrightarrow$	<i>no return value</i>	procedure

The `.../env` variants take an environment specified as a string→string alist. An environment of `#t` is taken to mean the current process' environment (*i.e.*, the value of the external char `**environ`).

[Rationale: `#f` is a more convenient marker for the current environment than `#t`, but would cause an ambiguity on Schemes that identify `#f` and `()`.]

The path-searching variants search the directories in the list `exec-path-list` for the program. A path-search is not performed if the program name contains a slash character—it is used directly. So a program with a name like `"bin/prog"` always executes the program `bin/prog` in the current working directory. See `$path` and `exec-path-list`, below.

Note that there is no analog to the C function `execv()`. To get the effect just do

```
(apply exec prog arglist)
```

All of these procedures flush buffered output and close unrevealed ports before executing the new binary. To avoid flushing buffered output, see `%exec` below.

Note that the C `exec()` procedure allows the zeroth element of the argument vector to be different from the file being executed, *e.g.*

```
char *argv[] = {"-", "-f", 0};
exec("/bin/csh", argv, envp);
```

The `scsh` `exec`, `exec-path`, `exec/env`, and `exec-path/env` procedures do not give this functionality—element 0 of the `arg` vector is always identical to the `prog` argument. In the rare case the user wishes to differentiate these two items, he can use the low-level `%exec` and `exec-path-search` procedures. These procedures never return under any circumstances. As with any other system call, if there is an error, they raise an exception.

`(%exec prog arglist env)`  $\longrightarrow$  *undefined* procedure  
`(exec-path-search fname pathlist)`  $\longrightarrow$  *string or #f* procedure

The `%exec` procedure is the low-level interface to the system call. The *arglist* parameter is a list of arguments; *env* is either a string $\rightarrow$ string alist or `#t`. The new program's `argv[0]` will be taken from `(car arglist)`, *not* from *prog*. An environment of `#t` means the current process' environment. `%exec` does not flush buffered output (see `flush-all-ports`).

All `exec` procedures, including `%exec`, coerce the *prog* and *arg* values to strings using the usual conversion rules: numbers are converted to decimal numerals, and symbols converted to their print-names.

`exec-path-search` searches the directories of *pathlist* looking for an occurrence of file *fname*. If no executable file is found, it returns `#f`. If *fname* contains a slash character, the path search is short-circuited, but the procedure still checks to ensure that the file exists and is executable—if not, it still returns `#f`. Users of this procedure should be aware that it invites a potential race condition: between checking the file with `exec-path-search` and executing it with `%exec`, the file's status might change. The only atomic way to do the search is to loop over the candidate file names, `exec`'ing each one and looping when the `exec` operation fails.

See `$path` and `exec-path-list`, below.

`(exit [status])`  $\longrightarrow$  *no return value* procedure  
`(%exit [status])`  $\longrightarrow$  *no return value* procedure

These procedures terminate the current process with a given exit status. The default exit status is 0. The low-level `%exit` procedure immediately terminates the process without flushing buffered output.

`(call-terminally thunk)`  $\longrightarrow$  *no return value* procedure

`call-terminally` calls its *thunk*. When the *thunk* returns, the process exits. Although `call-terminally` could be implemented as

`( $\lambda$  (thunk) (thunk) (exit 0))`

an implementation can take advantage of the fact that this procedure never returns. For example, the runtime can start with a fresh stack and

also start with a fresh dynamic environment, where shadowed bindings are discarded. This can allow the old stack and dynamic environment to be collected (assuming this data is not reachable through some live continuation).

(suspend)  $\longrightarrow$  *undefined* procedure  
 Suspend the current process with a SIGSTOP signal.

(fork [*thunk* or #f] [*continue-threads?*])  $\longrightarrow$  *proc* or #f procedure  
 (%fork [*thunk* or #f] [*continue-threads?*])  $\longrightarrow$  *proc* or #f procedure

fork with no arguments or #f instead of a thunk is like C fork(). In the parent process, it returns the child's *process object* (see below for more information on process objects). In the child process, it returns #f.

fork with an argument only returns in the parent process, returning the child's process object. The child process calls *thunk* and then exits.

fork flushes buffered output before forking, and sets the child process to non-interactive. %fork does not perform this bookkeeping; it simply forks.

The optional boolean argument *continue-threads?* specifies whether the currently active threads continue to run in the child or not. The default is #f.

(fork/pipe [*thunk*] [*continue-threads?*])  $\longrightarrow$  *proc* or #f procedure  
 (%fork/pipe [*thunk*] [*continue-threads?*])  $\longrightarrow$  *proc* or #f procedure

Like fork and %fork, but the parent and child communicate via a pipe connecting the parent's stdin to the child's stdout. These procedures side-effect the parent by changing his stdin.

In effect, fork/pipe splices a process into the data stream immediately upstream of the current process. This is the basic function for creating pipelines. Long pipelines are built by performing a sequence of fork/pipe calls. For example, to create a background two-process pipe a | b, we write:

```
(fork (λ () (fork/pipe a) (b)))
```

which returns the process object for b's process.

To create a background three-process pipe a | b | c, we write:

```
(fork (λ () (fork/pipe a)
              (fork/pipe b)
              (c)))
```



which returns the process object for *c*'s process.

Note that these procedures affect file descriptors, not ports. That is, the pipe is allocated connecting the child's file descriptor 1 to the parent's file descriptor 0. *Any previous Scheme port built over these affected file descriptors is shifted to a new, unused file descriptor with dup before allocating the I/O pipe.* This means, for example, that the ports bound to (current-input-port) and (current-output-port) in either process are not affected—they still refer to the same I/O sources and sinks as before. Remember the simple scsh rule: Scheme ports are bound to I/O sources and sinks, *not* particular file descriptors.

If the child process wishes to rebind the current output port to the pipe on file descriptor 1, it can do this using with-current-output-port or a related form. Similarly, if the parent wishes to change the current input port to the pipe on file descriptor 0, it can do this using set-current-input-port! or a related form. Here is an example showing how to set up the I/O ports on both sides of the pipe:

```
(fork/pipe (lambda ()
              (with-current-output-port (fdes->outport 1)
                (display "Hello, world.\n"))))

(set-current-input-port! (fdes->inport 0))
(read-line)           ; Read the string output by the child.
```

None of this is necessary when the I/O is performed by an exec'd program in the child or parent process, only when the pipe will be referenced by Scheme code through one of the default current I/O ports.

```
(fork/pipe+ conns [thunk] [continue-threads?]) → proc or #f procedure
(%fork/pipe+ conns [thunk] [continue-threads?]) → proc or #f procedure
```

Like fork/pipe, but the pipe connections between the child and parent are specified by the connection list *conns*. See the

```
(|+ conns pf1 ... pfn)
```

process form for a description of connection lists.

### 3.4.1 Process objects and process reaping

Scsh uses *process objects* to represent Unix processes. They are created by the fork procedure, and have the following exposed structure:

```
(define-record proc
  pid)
```

The only exposed slot in a proc record is the process' pid, the integer id assigned by Unix to the process. The only exported primitive procedures for manipulating process objects are `proc?` and `proc:pid`. Process objects are created with the `fork` procedure.

`(pid->proc pid [probe?])`  $\longrightarrow$  `proc` procedure

This procedure maps integer Unix process ids to scsh process objects. It is intended for use in interactive and debugging code, and is deprecated for use in production code. If there is no process object in the system indexed by the given pid, `pid->proc`'s action is determined by the *probe?* parameter (default `#f`):

<i>probe?</i>	Return
<code>#f</code>	<i>signal error condition.</i>
<code>'create</code>	Create new proc object.
True value	<code>#f</code>

Sometime after a child process terminates, scsh will perform a `wait` system call on the child in background, caching the process' exit status in the child's proc object. This is called "reaping" the process. Once the child has been waited, the Unix kernel can free the storage allocated for the dead process' exit information, so process reaping prevents the process table from becoming cluttered with un-waited dead child processes (a.k.a. "zombies"). This can be especially severe if the scsh process never waits on child processes at all; if the process table overflows with forgotten zombies, the OS may be unable to fork further processes.

Reaping a child process moves its exit status information from the kernel into the scsh process, where it is cached inside the child's process object. If the scsh user drops all pointers to the process object, it will simply be garbage collected. On the other hand, if the scsh program retains a pointer to the process object, it can use scsh's `wait` system call to synchronise with the child and retrieve its exit status multiple times (this is not possible with simple Unix integer pids in C—the programmer can only wait on a pid once).

Thus, process objects allow scsh programmer to do two things not allowed in other programming environments:

- Subprocesses that are never waited on are still removed from the process table, and their associated exit status data is eventually automatically garbage collected.
- Subprocesses can be waited on multiple times.

However, note that once a child has exited, if the scsh programmer drops all pointers to the child's proc object, the child's exit status will be reaped and

thrown away. This is the intended behaviour, and it means that integer pids are not enough to cause a process's exit status to be retained by the scsh runtime. (This is because it is clearly impossible to GC data referenced by integers.)

As a convenience for interactive use and debugging, all procedures that take process objects will also accept integer Unix pids as arguments, coercing them to the corresponding process objects. Since integer process ids are not reliable ways to keep a child's exit status from being reaped and garbage collected, programmers are encouraged to use process objects in production code.

`(autoreap-policy [policy])`  $\longrightarrow$  *old-policy* procedure

The scsh programmer can choose different policies for automatic process reaping. The policy is determined by applying this procedure to one of the values `'early`, `'late`, or `#f` (*i.e.*, no autoreap).

**early** The child is reaped from the Unix kernel's process table into scsh as soon as it dies. This is done by having a signal handler for the SIGCHLD signal reap the process.

**late** The child is not autoreaped until it dies *and* the scsh program drops all pointers to its process object. That is, the process table is cleaned out during garbage collection.

**#f** If autoreaping is turned off, process reaping is completely under control of the programmer, who can force outstanding zombies to be reaped by manually calling the `reap-zombies` procedure (see below).

Note that under any of the autoreap policies, a particular process *p* can be manually reaped into scsh by simply calling `(wait p)`. *All* zombies can be manually reaped with `reap-zombies`.

The `autoreap-policy` procedure returns the policy's previous value. Calling `autoreap-policy` with no arguments returns the current policy without no change.

`(reap-zombies)`  $\longrightarrow$  *boolean* procedure

This procedure reaps all outstanding exited child processes into scsh. It returns true if there are no more child processes to wait on, and false if there are outstanding processes still running or suspended.

### Issues with process reaping

Reaping a process does not reveal its process group at the time of death; this information is lost when the process reaped. This means that a dead, reaped

process is *not eligible* as a return value for a future `wait-process-group` call. This is not likely to be a problem for most code, as programs almost never wait on exited processes by process group. Process group waiting is usually applied to *stopped* processes, which are never reaped. So it is unlikely that this will be a problem for most programs.

Automatic process reaping is a useful programming convenience. However, if a program is careful to wait for all children, and does not wish automatic reaping to happen, the programmer can simply turn process autoreaping off.

Programs that do not wish to use automatic process reaping should be aware that some `scsh` routines create subprocesses but do not return the child's pid: `run/port*`, and its related procedures and special forms (`run/strings`, *et al.*). Automatic process reaping will clean the child processes created by these procedures out of the kernel's process table. If a program doesn't use process reaping, it should either avoid these forms, or use `wait-any` to wait for the children to exit.

### 3.4.2 Process waiting

`(wait proc/pid [flags])`  $\longrightarrow$  *status* procedure

This procedure waits until a child process exits, and returns its exit code. The *proc/pid* argument is either a process object (section 3.4.1) or an integer process id. `wait` returns the child's exit status code (or suspension code, if the `wait/stopped-children` option is used, see below). Status values can be queried with the procedures in section 3.4.3.

The *flags* argument is an integer whose bits specify additional options. It is composed by or'ing together the following flags:

Flag	Meaning
<code>wait/poll</code>	Return <code>#f</code> immediately if child still active.
<code>wait/stopped-children</code>	Wait for suspend as well as exit.

`(wait-any [flags])`  $\longrightarrow$  [*proc status*] procedure

The optional *flags* argument is as for `wait`. This procedure waits for any child process to exit (or stop, if the `wait/stopped-children` flag is used). It returns the process' process object and status code. If there are no children left for which to wait, the two values `[#f #t]` are returned. If the `wait/poll` flag is used, and none of the children are immediately eligible for waiting, then the values `[#f #f]` are returned:

<code>[#f #f]</code>	Poll, none ready
<code>[#f #t]</code>	No children

Wait-any will not return a process that has been previously waited by any other process-wait procedure (wait, wait-any, and wait-process-group). It will return reaped processes that haven't yet been waited.

The use of wait-any is deprecated.

(wait-process-group *proc/pid* [*flags*]) → [*proc status*] procedure

This procedure waits for any child whose process group is *proc/pid* (either a process object or a pid). The *flags* argument is as for wait.

Note that if the programmer wishes to wait for exited processes by process group, the program should take care not to use process reaping (section 3.4.1), as this loses process group information. However, most process-group waiting is for stopped processes (to implement job control), so this is rarely an issue, as stopped processes are not subject to reaping.

### 3.4.3 Analysing process status codes

When a child process dies (or is suspended), its parent can call the wait procedure to recover the exit (or suspension) status of the child. The exit status is a small integer that encodes information describing how the child terminated. The bit-level format of the exit status is not defined by POSIX; you must use the following three functions to decode one. However, if a child terminates normally with exit code 0, POSIX does require wait to return an exit status that is exactly zero. So (zero? *status*) is a correct way to test for non-error, normal termination, *e.g.*,

```
(if (zero? (run (rcp scsh.tar.gz lambda.csd.hku.hk:)))
    (delete-file "scsh.tar.gz"))
```

(status:exit-val *status*) → *integer or #f* procedure

(status:stop-sig *status*) → *integer or #f* procedure

(status:term-sig *status*) → *integer or #f* procedure

For a given status value produced by calling wait, exactly one of these routines will return a true value.

If the child process exited normally, status:exit-val returns the exit code for the child process (*i.e.*, the value the child passed to exit or returned from main). Otherwise, this function returns false.

If the child process was suspended by a signal, `status:stop-sig` returns the signal that suspended the child. Otherwise, this function returns `false`.

If the child process terminated abnormally, `status:term-sig` returns the signal that terminated the child. Otherwise, this function returns `false`.

### 3.5 Process state

<code>(umask)</code>	$\longrightarrow$ <i>fixnum</i>	procedure
<code>(set-umask perms)</code>	$\longrightarrow$ <i>undefined</i>	procedure
<code>(with-umask* perms thunk)</code>	$\longrightarrow$ <i>value(s) of thunk</i>	procedure
<code>(with-umask perms . body)</code>	$\longrightarrow$ <i>value(s) of body</i>	syntax

The process' current umask is retrieved with `umask`, and set with `(set-umask perms)`. Calling `with-umask*` changes the umask to *perms* for the duration of the call to *thunk*. If the program throws out of *thunk* by invoking a continuation, the umask is reset to its external value. If the program throws back into *thunk* by calling a stored continuation, the umask is restored to the *perms* value. The special form `with-umask` is equivalent in effect to the procedure `with-umask*`, but does not require the programmer to explicitly wrap a `( $\lambda$  () ...)` around the body of the code to be executed.

<code>(chdir [fname])</code>	$\longrightarrow$ <i>undefined</i>	procedure
<code>(cwd)</code>	$\longrightarrow$ <i>string</i>	procedure
<code>(with-cwd* fname thunk)</code>	$\longrightarrow$ <i>value(s) of thunk</i>	procedure
<code>(with-cwd fname . body)</code>	$\longrightarrow$ <i>value(s) of body</i>	syntax

These forms manipulate the current working directory. The `cwd` can be changed with `chdir` (although in most cases, `with-cwd` is preferable). If `chdir` is called with no arguments, it changes the `cwd` to the user's home directory. The `with-cwd*` procedure calls *thunk* with the `cwd` temporarily set to *fname*; when *thunk* returns, or is exited in a non-local fashion (e.g., by raising an exception or by invoking a continuation), the `cwd` is returned to its original value. The special form `with-cwd` is simply syntactic sugar for `with-cwd*`.

<code>(pid)</code>	$\longrightarrow$ <i>fixnum</i>	procedure
<code>(parent-pid)</code>	$\longrightarrow$ <i>fixnum</i>	procedure
<code>(process-group)</code>	$\longrightarrow$ <i>fixnum</i>	procedure
<code>(set-process-group [proc/pid] pgrp)</code>	$\longrightarrow$ <i>undefined</i>	procedure

`(pid)` and `(parent-pid)` retrieve the process id for the current process and its parent. `(process-group)` returns the process group of the current

process. A process' process-group can be set with `set-process-group`; the value *proc/pid* specifies the affected process. It may be either a process object or an integer process id, and defaults to the current process.

<code>(set-priority which who priority)</code>	$\longrightarrow$ <i>undefined</i>	procedure
<code>(priority which who)</code>	$\longrightarrow$ <i>fixnum</i>	procedure
<code>(nice [proc/pid delta])</code>	$\longrightarrow$ <i>undefined</i>	procedure

These procedures set and access the priority of processes. I can't remember how `set-priority` and `priority` work, so no documentation, and besides, they aren't implemented yet, anyway.

<code>(user-login-name)</code>	$\longrightarrow$ <i>string</i>	procedure
<code>(user-uid)</code>	$\longrightarrow$ <i>fixnum</i>	procedure
<code>(user-gid)</code>	$\longrightarrow$ <i>fixnum</i>	procedure
<code>(user-supplementary-gids)</code>	$\longrightarrow$ <i>fixnum list</i>	procedure
<code>(set-uid uid)</code>	$\longrightarrow$ <i>undefined</i>	procedure
<code>(set-gid gid)</code>	$\longrightarrow$ <i>undefined</i>	procedure

These routines get and set the effective and real user and group ids. The `set-uid` and `set-gid` routines correspond to the POSIX `setuid()` and `setgid()` procedures.

<code>(user-effective-uid)</code>	$\longrightarrow$ <i>fixnum</i>	procedure
<code>(set-user-effective-uid fixnum)</code>	$\longrightarrow$ <i>undefined</i>	procedure
<code>(with-user-effective-uid* fixnum thunk)</code>	$\longrightarrow$ <i>value(s) of thunk</i>	procedure
<code>(with-user-effective-uid fixnum . body)</code>	$\longrightarrow$ <i>value(s) of body</i>	syntax
<code>(user-effective-gid)</code>	$\longrightarrow$ <i>fixnum</i>	procedure
<code>(set-user-effective-gid fixnum)</code>	$\longrightarrow$ <i>undefined</i>	procedure
<code>(with-user-effective-gid* fixnum thunk)</code>	$\longrightarrow$ <i>value(s) of thunk</i>	procedure
<code>(with-user-effective-gid fixnum . body)</code>	$\longrightarrow$ <i>value(s) of body</i>	syntax

These forms manipulate the effective user/group IDs. Possible values for setting this resource are either the real user/group ID or the saved set-user/group-ID. The `with-...` forms perform the usual temporary assignment during the execution of the second argument. The effective user and group IDs are thread-local.

<code>(process-times)</code>	$\longrightarrow$ [ <i>fixnum fixnum fixnum fixnum</i> ]	procedure
------------------------------	--	-----------

Returns four values:

- user CPU time in clock-ticks
- system CPU time in clock-ticks
- user CPU time of all descendant processes
- system CPU time of all descendant processes

Note that CPU time clock resolution is not the same as the real-time clock resolution provided by `time+ticks`. That's Unix.

(cpu-ticks/sec)  $\longrightarrow$  *integer* procedure

Returns the resolution of the CPU timer in clock ticks per second. This can be used to convert the times reported by `process-times` to seconds.

### 3.6 User and group database access

These procedures are used to access the user and group databases (*e.g.*, the ones traditionally stored in `/etc/passwd` and `/etc/group`.)

(user-info *uid/name*)  $\longrightarrow$  *record* procedure

Return a `user-info` record giving the recorded information for a particular user:

```
(define-record user-info
  name uid gid home-dir shell)
```

The *uid/name* argument is either an integer `uid` or a string `user-name`.

(->uid *uid/name*)  $\longrightarrow$  *fixnum* procedure

(->username *uid/name*)  $\longrightarrow$  *string* procedure

These two procedures coerce integer `uid`'s and user names to a particular form.

(group-info *gid/name*)  $\longrightarrow$  *record* procedure

Return a `group-info` record giving the recorded information for a particular group:

```
(define-record group-info
  name gid members)
```

The *gid/name* argument is either an integer `gid` or a string `group-name`.

### 3.7 Accessing command-line arguments

`command-line-arguments` *string list*

(command-line)  $\longrightarrow$  *string list* procedure

The list of strings `command-line-arguments` contains the arguments passed to the `scsh` process on the command line. Calling `(command-line)` returns the complete `argv` string list, including the program. So if we run a `scsh` program

```
/usr/shivers/bin/myls -CF src
```



then `command-line-arguments` is

```
("-CF" "src")
```

and `(command-line)` returns

```
("/usr/shivers/bin/myls" "-CF" "src")
```

`command-line` returns a fresh list each time it is called. In this way, the programmer can get a fresh copy of the original argument list if `command-line-arguments` has been modified or is lexically shadowed.

<code>(arg arglist n [default])</code>	$\longrightarrow$ <i>string</i>	procedure
<code>(arg* arglist n [default-thunk])</code>	$\longrightarrow$ <i>string</i>	procedure
<code>(argv n [default])</code>	$\longrightarrow$ <i>string</i>	procedure

These procedures are useful for accessing arguments from argument lists. `arg` returns the  $n^{\text{th}}$  element of *arglist*. The index is 1-based. If  $n$  is too large, *default* is returned; if no *default*, then an error is signaled.

`arg*` is similar, except that the *default-thunk* is called to generate the default value.

`(argv n)` is simply `(arg (command-line) (+ n 1))`. The +1 offset ensures that the two forms

```
(arg command-line-arguments n)
(argv n)
```

return the same argument (assuming the user has not rebound or modified `command-line-arguments`).

Example:

```
(if (null? command-line-arguments)
    (& (xterm -n ,host -title ,host
           -name ,(string-append "xterm_" host)))
    (let* ((progname (file-name-nondirectory (argv 1)))
           (title (string-append host ":" progname)))
      (& (xterm -n      ,title
                -title ,title
                -e      ,@command-line-arguments))))
```

A subtlety: when the `scsh` interpreter is used to execute a `scsh` program, the program name reported in the head of the `(command-line)` list is the `scsh` program, *not* the interpreter. For example, if we have a shell script in file `fullecho`:

```
#!/usr/local/bin/scsh -s
!#
(for-each (λ (arg) (display arg) (display " "))
  (command-line))
```

and we run the program

```
fullecho hello world
```

the program will print out

```
fullecho hello world
```

not

```
/usr/local/bin/scsh -s fullecho hello world
```

This argument line processing ensures that if a scsh program is subsequently compiled into a standalone executable or byte-compiled to a heap-image executable by the Scheme 48 virtual machine, its semantics will be unchanged—the arglist processing is invariant. In effect, the

```
/usr/local/bin/scsh -s
```

is not part of the program; it’s a specification for the machine to execute the program on, so it is not properly part of the program’s argument list.

### 3.8 System parameters

(system-name)  $\longrightarrow$  *string* procedure

Returns the name of the host on which we are executing. This may be a local name, such as “solar,” as opposed to a fully-qualified domain name such as “solar.csie.ntu.edu.tw.”

(uname)  $\longrightarrow$  *uname-record* procedure

Returns a *uname-record* of the following structure:

```
(define-record uname
  os-name
  node-name
  release
  version
  machine)
```

Each of the fields contains a string.

Be aware that POSIX limits the length of all entries to 32 characters, and that the node name does not necessarily correspond to the fully-qualified domain name.

### 3.9 Signal system

Signal numbers are bound to the variables `signal/hup`, `signal/int`, .... See tables 3.9 and 3.3 for the full list.

`(signal-process proc sig)`  $\longrightarrow$  *undefined* procedure  
`(signal-process-group prgrp sig)`  $\longrightarrow$  *undefined* procedure

These two procedures send signals to a specific process, and all the processes in a specific process group, respectively. The *proc* and *prgrp* arguments are either processes or integer process ids.

`(itimer secs)`  $\longrightarrow$  *undefined* procedure

Schedules a timer interrupt in *secs* seconds.

{Note A}s the thread system needs the timer interrupt for its own purpose, `itimer` works by spawning a thread which calls the interrupt handler for `interrupt/alm` after the specified time.

`(process-sleep secs)`  $\longrightarrow$  *undefined* procedure  
`(process-sleep-until time)`  $\longrightarrow$  *undefined* procedure

The `sleep` procedure causes the process to sleep for *secs* seconds. The `sleep-until` procedure causes the process to sleep until *time* (see section 3.10).

{Note The use of these procedures is deprecated as they suspend *all* running threads, including the ones Scsh uses for administrative purposes. Consider using the `sleep` procedure from the `thread` package.}

#### Interrupt handlers

Scsh interrupt handlers are complicated by the fact that `scsh` is implemented on top of the Scheme 48 virtual machine, which has its own interrupt system, independent of the Unix signal system. This means that Unix signals are delivered in two stages: first, Unix delivers the signal to the Scheme 48 virtual machine, then the Scheme 48 virtual machine delivers the signal to the executing Scheme program as a Scheme 48 interrupt. This ensures that signal delivery happens between two vm instructions, keeping individual instructions atomic.

The Scheme 48 machine has its own set of interrupts, which includes the asynchronous Unix signals (table 3.9). Note that `scsh` does *not* support signal handlers for “synchronous” Unix signals, such as `signal/ill` or `signal/pipe` (see table 3.3). Synchronous occurrences of these signals are better handled by raising a Scheme exception. We recommend you avoid using signal handlers unless you absolutely have to; Section 9.4 describes a better interface to signals.

`(signal->interrupt integer)`  $\longrightarrow$  *integer* procedure

Interrupt	Unix signal	OS Variant
interrupt/alm <sup>a</sup>	signal/alm	POSIX
interrupt/int <sup>b</sup>	signal/int	POSIX
interrupt/memory-shortage	N/A	
interrupt/chld	signal/chld	POSIX
interrupt/cont	signal/cont	POSIX
interrupt/hup	signal/hup	POSIX
interrupt/quit	signal/quit	POSIX
interrupt/term	signal/term	POSIX
interrupt/tstp	signal/tstp	POSIX
interrupt/usr1	signal/usr1	POSIX
interrupt/usr2	signal/usr2	POSIX
interrupt/info	signal/info	BSD only
interrupt/io	signal/io	BSD + SVR4
interrupt/poll	signal/poll	SVR4 only
interrupt/prof	signal/prof	BSD + SVR4
interrupt/pwr	signal/pwr	SVR4 only
interrupt/urg	signal/urg	BSD + SVR4
interrupt/vtalrm	signal/vtalrm	BSD + SVR4
interrupt/winch	signal/winch	BSD + SVR4
interrupt/xcpu	signal/xcpu	BSD + SVR4
interrupt/xfsz	signal/xfsz	BSD + SVR4

Table 3.2: Scheme 48 virtual-machine interrupts and related Unix signals. Only the POSIX signals are guaranteed to be defined; however, your implementation and OS may define other signals and interrupts not listed here.

<sup>a</sup> Also bound to Scheme 48 interrupt `interrupt/alarm`.

<sup>b</sup> Also bound to Scheme 48 interrupt `interrupt/keyboard`.

Unix signal	Type	OS Variant
signal/stop	Uncatchable	POSIX
signal/kill	Uncatchable	POSIX
signal/abrt	Synchronous	POSIX
signal/fpe	Synchronous	POSIX
signal/ill	Synchronous	POSIX
signal/pipe	Synchronous	POSIX
signal/segv	Synchronous	POSIX
signal/ttin	Synchronous	POSIX
signal/ttou	Synchronous	POSIX
signal/bus	Synchronous	BSD + SVR4
signal/emt	Synchronous	BSD + SVR4
signal/iot	Synchronous	BSD + SVR4
signal/sys	Synchronous	BSD + SVR4
signal/trap	Synchronous	BSD + SVR4

Table 3.3: Uncatchable and synchronous Unix signals. While these signals may be sent with `signal-process` or `signal-process-group`, there are no corresponding `scsh` interrupt handlers. Only the POSIX signals are guaranteed to be defined; however, your implementation and OS may define other signals not listed here.

The programmer maps from Unix signals to Scheme 48 interrupts with the `signal->interrupt` procedure. If the signal does not have a defined Scheme 48 interrupt, an error is signaled.

`(interrupt-set integer1 ... integern)`  $\longrightarrow$  *integer* procedure

This procedure builds interrupt sets from its interrupt arguments. A set is represented as an integer using a two's-complement representation of the bit set.

`(enabled-interrupts)`  $\longrightarrow$  *interrupt-set* procedure

`(set-enabled-interrupts interrupt-set)`  $\longrightarrow$  *interrupt-set* procedure

Get and set the value of the enabled-interrupt set. Only interrupts in this set have their handlers called when delivered. When a disabled interrupt is delivered to the Scheme 48 machine, it is held pending until it becomes enabled, at which time its handler is invoked.

Interrupt sets are represented as integer bit sets (constructed with the `interrupt-set` function). The `set-enabled-interrupts` procedure returns the previous value of the enabled-interrupt set.

`(with-enabled-interrupts interrupt-set . body)`  $\longrightarrow$  *value(s) of body* syntax

`(with-enabled-interrupts* interrupt-set thunk)`  $\longrightarrow$  *value(s) of thunk* procedure

Run code with a given set of interrupts enabled. Note that “enabling” an interrupt means enabling delivery from the Scheme 48 vm to the scsh program. Using the Scheme 48 interrupt system is fairly lightweight, and does not involve actually making a system call. Note that enabling an interrupt means that the assigned interrupt handler is allowed to run when the interrupt is delivered. Interrupts not enabled are held pending when delivered.

Interrupt sets are represented as integer bit sets (constructed with the `interrupt-set` function).

`(set-interrupt-handler interrupt handler)`  $\longrightarrow$  *old-handler* procedure

Assigns a handler for a given interrupt, and returns the interrupt's old handler. The *handler* argument is `#f` (ignore), `#t` (default), or a procedure taking an integer argument; the return value follows the same conventions. Note that the *interrupt* argument is an interrupt value, not a signal value. An interrupt is delivered to the Scheme 48 machine by (1) blocking all interrupts, and (2) applying the handler procedure to the set of interrupts that were enabled prior to the interrupt delivery. If the procedure returns normally (*i.e.*, it doesn't throw to a continuation), the set

of enabled interrupts will be returned to its previous value. (To restore the enabled-interrupt set before throwing out of an interrupt handler, see `set-enabled-interrupts`)

{Note If you set a handler for the `interrupt/chld` interrupt, you may break `scsh`'s autoreaping process machinery. See the discussion of autoreaping in section 3.4.1.}

`(interrupt-handler interrupt)`  $\longrightarrow$  *handler* procedure

Return the handler for a given interrupt. Note that the argument is an interrupt value, not a signal value. A handler is either `#f` (ignore), `#t` (default), or a procedure taking an integer argument.

## 3.10 Time

`Scsh`'s time system is fairly sophisticated, particularly with respect to its careful treatment of time zones. However, casual users shouldn't be intimidated; all of the complexity is optional, and defaulting all the optional arguments reduces the system to a simple interface.

### 3.10.1 Terminology

"UTC" and "UCT" stand for "universal coordinated time," which is the official name for what is colloquially referred to as "Greenwich Mean Time."

POSIX allows a single time zone to specify *two* different offsets from UTC: one standard one, and one for "summer time." Summer time is frequently some sort of daylight savings time.

The `scsh` time package consistently uses this terminology: we never say "gmt" or "dst;" we always say "utc" and "summer time."

### 3.10.2 Basic data types

We have two types: *time* and *date*.

A *time* specifies an instant in the history of the universe. It is location and time-zone independent.<sup>4</sup> A time is a real value giving the number of elapsed seconds since the Unix "epoch" (Midnight, January 1, 1970 UTC). Time values

---

<sup>4</sup>Physics pedants please note: The `scsh` authors live in a Newtonian universe. We disclaim responsibility for calculations performed in non-ANSI standard light-cones.

provide arbitrary time resolution, limited only by the number system of the underlying Scheme system.

A *date* is a name for an instant in time that is specified relative to some location/time-zone in the world, *e.g.*:

Friday October 31, 1994 3:47:21 pm EST.

Dates provide one-second resolution, and are expressed with the following record type:

```
(define-record date      ; A Posix tm struct
  seconds      ; Seconds after the minute [0-59]
  minute       ; Minutes after the hour [0-59]
  hour         ; Hours since midnight [0-23]
  month-day    ; Day of the month [1-31]
  month        ; Months since January [0-11]
  year         ; Years since 1900
  tz-name      ; Time-zone name: #f or a string.
  tz-secs      ; Time-zone offset: #f or an integer.
  summer?      ; Summer (Daylight Savings) time in effect?
  week-day     ; Days since Sunday [0-6]
  year-day     ; Days since Jan. 1 [0-365])
```

If the `tz-secs` field is given, it specifies the time-zone's offset from UTC in seconds. If it is specified, the `tz-name` and `summer?` fields are ignored when using the date structure to determine a specific instant in time.

If the `tz-name` field is given, it is a time-zone string such as "EST" or "HKT" understood by the OS. Since POSIX time-zone strings can specify dual standard/summer time-zones (e.g., "EST5EDT" specifies U.S. Eastern Standard/Eastern Daylight Time), the value of the `summer?` field is used to resolve the ambiguous boundary cases. For example, on the morning of the Fall daylight savings change-over, 1:00am–2:00am happens twice. Hence the date 1:30 am on this morning can specify two different seconds; the `summer?` flag says which one.

A date with `tz-name = tz-secs = #f` is a date that is specified in terms of the system's current time zone.

There is redundancy in the date data structure. For example, the `year-day` field is redundant with the `month-day` and `month` fields. Either of these implies the values of the `week-day` field. The `summer?` and `tz-name` fields are redundant with the `tz-secs` field in terms of specifying an instant in time. This redundancy is provided because consumers of dates may want it broken out in different ways. The `scsh` procedures that produce date records fill them out completely. However, when date records produced by the programmer are passed to `scsh` procedures, the redundancy is resolved by ignoring some of the secondary fields. This is described for each procedure below.



`(make-date s min h mday mon y [tzn tzs summ? wday yday])`  $\longrightarrow$  *date* procedure

When making a `date` record, the last five elements of the record are optional, and default to `#f`, `#f`, `#f`, `0`, and `0` respectively. This is useful when creating a `date` record to pass as an argument to `time`. Other procedures, however, may refuse to work with these incomplete `date` records.

### 3.10.3 Time zones

Several time procedures take time zones as arguments. When optional, the time zone defaults to local time zone. Otherwise the time zone can be one of:

<code>#f</code>	Local time
<code>Integer</code>	Seconds of offset from UTC. For example, New York City is -18000 (-5 hours), San Francisco is -28800 (-8 hours).
<code>String</code>	A POSIX time zone string understood by the OS ( <i>i.e.</i> , the sort of time zone assigned to the <code>\$TZ</code> environment variable).

An integer time zone gives the number of seconds you must add to UTC to get time in that zone. It is *not* “seconds west” of UTC—that flips the sign.

To get UTC time, use a time zone of either `0` or `"UCT0"`.

### 3.10.4 Procedures

`(time+ticks)`  $\longrightarrow$  *[secs ticks]* procedure  
`(ticks/sec)`  $\longrightarrow$  *real* procedure

The current time, with sub-second resolution. Sub-second resolution is not provided by POSIX, but is available on many systems. The time is returned as elapsed seconds since the Unix epoch, plus a number of sub-second “ticks.” The length of a tick may vary from implementation to implementation; it can be determined from `(ticks/sec)`.

The system clock is not required to report time at the full resolution given by `(ticks/sec)`. For example, on BSD, time is reported at  $1\mu\text{s}$  resolution, so `(ticks/sec)` is 1,000,000. That doesn’t mean the system clock has micro-second resolution.

If the OS does not support sub-second resolution, the `ticks` value is always 0, and `(ticks/sec)` returns 1.

*Remark:* I chose to represent system clock resolution as `ticks/sec` instead of `sec/tick` to increase the odds that the value could be represented as an exact integer, increasing efficiency and making it easier

for Scheme implementations that don't have sophisticated numeric support to deal with the quantity.

You can convert seconds and ticks to seconds with the expression

```
(+ secs (/ ticks (ticks/sec)))
```

Given that, why not have the fine-grain time procedure just return a non-integer real for time? Following Common Lisp, I chose to allow the system clock to report sub-second time in its own units to lower the overhead of determining the time. This would be important for a system that wanted to precisely time the duration of some event. Time stamps could be collected with little overhead, deferring the overhead of precisely calculating with them until after collection.

This is all a bit academic for the Scheme 48 implementation, where we determine time with a heavyweight system call, but it's nice to plan for the future.

(date)	→	<i>date-record</i>	procedure
(date [time tz])	→	<i>date-record</i>	procedure

Simple (date) returns the current date, in the local time zone.

With the optional arguments, date converts the time to the date as specified by the time zone *tz*. *Time* defaults to the current time; *tz* defaults to local time, and is as described in the time-zone section.

If the *tz* argument is an integer, the date's *tz-name* field is a POSIX time zone of the form "UTC+*hh*:*mm*:*ss*"; the trailing *mm*:*ss* portion is deleted if it is zeroes.

*Oops:* The Posix facility for converting dates to times, `mktime()`, has a broken design: it indicates an error by returning -1, which is also a legal return value (for date 23:59:59 UCT, 12/31/1969). Scsh resolves the ambiguity in a paranoid fashion: it always reports an error if the underlying Unix facility returns -1. We feel your pain.

(time)	→	<i>integer</i>	procedure
(time [date])	→	<i>integer</i>	procedure

Simple (time) returns the current time.

With the optional date argument, time converts a date to a time. *Date* defaults to the current date.

Note that the input *date* record is overconstrained. *time* ignores *date*'s week-day and year-day fields. If the date's *tz-secs* field is set, the *tz-name* and *summer?* fields are ignored.

#f	Resolve an ambiguous time in favor of non-summer time.
true	Resolve an ambiguous time in favor of summer time.

Outside of these boundary cases, the `summer?` flag is ignored. For example, if the standard/summer change-overs happen in the Fall and the Spring, then the value of `summer?` is ignored for a January or July date. A January date would be resolved with standard time, and a July date with summer time, regardless of the `summer?` value.

The `summer?` flag is also ignored if the time zone doesn't have a summer time—for example, simple UTC.

Date->string formats the date as a 24-character string of the form:

Format-date formats the date according to the format string *fmt*. The format string is copied verbatim, except that tilde characters indicate conversion specifiers that are replaced by fields from the date record. Figure 3.1 gives the full set of conversion specifiers supported by format-date.

This procedure fills in missing, redundant slots in a date record. In decreasing order of priority:

- If `tz-secs` is defined, but `tz-name` is not, it is assigned a time-zone name of the form “UTC+*hh:mm:ss*”; the trailing *:mm:ss* portion is deleted if it is zeroes.

~ ~    Converted to the ~ character.  
 ~a    abbreviated weekday name  
 ~A    full weekday name  
 ~b    abbreviated month name  
 ~B    full month name  
 ~c    time and date using the time and date representation for the locale  
       (~X ~x)  
 ~d    day of the month as a decimal number (01-31)  
 ~H    hour based on a 24-hour clock as a decimal number (00-23)  
 ~I    hour based on a 12-hour clock as a decimal number (01-12)  
 ~j    day of the year as a decimal number (001-366)  
 ~m    month as a decimal number (01-12)  
 ~M    minute as a decimal number (00-59)  
 ~p    AM/PM designation associated with a 12-hour clock  
 ~S    second as a decimal number (00-61)  
 ~U    week number of the year; Sunday is first day of week (00-53)  
 ~w    weekday as a decimal number (0-6), where Sunday is 0  
 ~W    week number of the year; Monday is first day of week (00-53)  
 ~x    date using the date representation for the locale  
 ~X    time using the time representation for the locale  
 ~y    year without century (00-99)  
 ~Y    year with century (*e.g.* 1990)  
 ~Z    time zone name or abbreviation, or no characters if no time zone is  
       determinable

Figure 3.1: format-date conversion specifiers

- **tz-name, date, summer?  $\Rightarrow$  tz-secs, summer?**

If the date information is provided up to second resolution, tz-name is also provided, and tz-secs is not set, then tz-secs and summer? are set to their correct values. Summer-time ambiguities are resolved using the original value of summer?. If the time zone doesn't have a summer time variant, then summer? is set to #f.

- **local time, date, summer?  $\Rightarrow$  tz-name, tz-secs, summer?**

If the date information is provided up to second resolution, but no time zone information is provided (both tz-name and tz-secs aren't set), then we proceed as in the above case, except the system's current time zone is used.

These rules allow one particular ambiguity to escape: if both tz-name and tz-secs are set, they are not brought into agreement. It isn't clear how to do this, nor is it clear which one should take precedence.

*Oops: fill-in-date! isn't implemented yet.*

### 3.11 Environment variables

(setenv *var val*)  $\longrightarrow$  *undefined* procedure  
 (getenv *var*)  $\longrightarrow$  *string* procedure

These functions get and set the process environment, stored in the external C variable `char **environ`. An environment variable *var* is a string. If an environment variable is set to a string *val*, then the process' global environment structure is altered with an entry of the form "*var=val*". If *val* is #f, then any entry for *var* is deleted.

(env->alist)  $\longrightarrow$  *string* $\rightarrow$ *string* *alist* procedure

The env->alist procedure converts the entire environment into an alist, e.g.,

```
(("TERM" . "vt100")
 ("SHELL" . "/usr/local/bin/scsh")
 ("PATH" . "/sbin:/usr/sbin:/bin:/usr/bin")
 ("EDITOR" . "emacs")
 ...)
```

(alist->env *alist*)  $\longrightarrow$  *undefined* procedure

*Alist* must be an alist whose keys are all strings, and whose values are all either strings or string lists. String lists are converted to colon lists (see below). The alist is installed as the current Unix environment (*i.e.*, converted to a null-terminated C vector of "*var=val*" strings which is assigned to the global char **\*\*environ**).

```
;;; Note $PATH entry is converted
;;; to /sbin:/usr/sbin:/bin:/usr/bin.
(alist->env '("TERM" . "vt100")
            ("PATH" "/sbin" "/usr/sbin" "/bin")
            ("SHELL" . "/usr/local/bin/scsh"))
```

Note that `env->alist` and `alist->env` are not exact inverses—`alist->env` will convert a list value into a single colon-separated string, but `env->alist` will not parse colon-separated values into lists. (See the `$PATH` element in the examples given for each procedure.)

The following three functions help the programmer manipulate alist tables in some generally useful ways. They are all defined using `equal?` for key comparison.

`(alist-delete key alist)`  $\longrightarrow$  *alist* procedure  
Delete any entry labelled by value *key*.

`(alist-update key val alist)`  $\longrightarrow$  *alist* procedure  
Delete *key* from *alist*, then cons on a (*key* . *val*) entry.

`(alist-compress alist)`  $\longrightarrow$  *alist* procedure  
Compresses *alist* by removing shadowed entries. Example:

```
;;; Shadowed (1 . c) entry removed.
(alist-compress '( (1 . a) (2 . b) (1 . c) (3 . d) ))
 $\Rightarrow$  ((1 . a) (2 . b) (3 . d))
```

`(with-env* env-alist-delta thunk)`  $\longrightarrow$  *value(s) of thunk* procedure  
`(with-total-env* env-alist thunk)`  $\longrightarrow$  *value(s) of thunk* procedure

These procedures call *thunk* in the context of an altered environment. They return whatever values *thunk* returns. Non-local returns restore the environment to its outer value; throwing back into the *thunk* by invoking a stored continuation restores the environment back to its inner value.

The *env-alist-delta* argument specifies a *modification* to the current environment—*think*'s environment is the original environment overridden with the bindings specified by the alist delta.

The *env-alist* argument specifies a complete environment that is installed for *think*.

```
(with-env env-alist-delta . body)  → value(s) of body      syntax
(with-total-env env-alist . body)  → value(s) of body      syntax
```

These special forms provide syntactic sugar for *with-env\** and *with-total-env\**. The env alists are not evaluated positions, but are implicitly backquoted. In this way, they tend to resemble binding lists for *let* and *let\** forms.

Example: These four pieces of code all run the mailer with special \$TERM and \$EDITOR values.

```
(with-env (("TERM" . "xterm") ("EDITOR" . ,my-editor))
  (run (mail shivers@lcs.mit.edu)))

(with-env* `(("TERM" . "xterm") ("EDITOR" . ,my-editor))
  (λ () (run (mail shivers@csd.hku.hk))))

(run (begin (setenv "TERM" "xterm")          ; Env mutation happens
            (setenv "EDITOR" my-editor) ; in the subshell.
            (exec-epf (mail shivers@research.att.com)))))

;; In this example, we compute an alternate environment ENV2
;; as an alist, and install it with an explicit call to the
;; EXEC-PATH/ENV procedure.
(let* ((env (env->alist))                ; Get the current environment,
      (env1 (alist-update env "TERM" "xterm")) ; and compute
      (env2 (alist-update env1 "EDITOR" my-editor))) ; the new env.
  (run (begin (exec-path/env "mail" env2 "shivers@cs.cmu.edu"))))
```

### 3.11.1 Path lists and colon lists

When environment variables such as \$PATH need to encode a list of strings (such as a list of directories to be searched), the common Unix convention is to separate the list elements with colon delimiters.<sup>5</sup> To convert between the colon-separated string encoding and the list-of-strings representation, see the *infix-splitter* function (section 8.1.2) and the string library's *string-join* function. For example,

<sup>5</sup>...and hope the individual list elements don't contain colons themselves.

```
(define split (infix-splitter (rx ":")))
(split "/sbin:/bin:./usr/bin") ⇒
'("/sbin" "/bin" "" "/usr/bin")
(string-join ":" '("/sbin" "/bin" "" "/usr/bin")) ⇒
"/sbin:/bin:./usr/bin"
```

The following two functions are useful for manipulating these ordered lists, once they have been parsed from their colon-separated form.

```
(add-before elt before list)  → list      procedure
(add-after  elt after list)   → list      procedure
```

These functions are for modifying search-path lists, where element order is significant.

`add-before` adds *elt* to the list immediately before the first occurrence of *before* in the list. If *before* is not in the list, *elt* is added to the end of the list.

`add-after` is similar: *elt* is added after the last occurrence of *after*. If *after* is not found, *elt* is added to the beginning of the list.

Neither function destructively alters the original path-list. The result may share structure with the original list. Both functions use `equal?` for comparing elements.

### 3.11.2 \$USER, \$HOME, and \$PATH

Like `sh` and unlike `csh`, `scsh` has *no* interactive dependencies on environment variables. It does, however, initialise certain internal values at startup time from the initial process environment, in particular `$HOME` and `$PATH`. `Scsh` never uses `$USER` at all. It computes (user-login-name) from the system call (user-uid).

```
home-directory      string
exec-path-list      string list thread-fluid
```

`Scsh` accesses `$HOME` at start-up time, and stores the value in the global variable `home-directory`. It uses this value for `~` lookups and for returning to home on (`chdir`).

`Scsh` accesses `$PATH` at start-up time, colon-splits the path list, and stores the value in the thread fluid `exec-path-list`. This list is used for `exec-path` and `exec-path/env` searches.

To access, rebind or side-effect thread-fluid cells, you must open the `thread-fluids` package.



## 3.12 Terminal device control

Scsh provides a complete set of routines for manipulating terminal devices—putting them in “raw” mode, changing and querying their special characters, modifying their I/O speeds, and so forth. The scsh interface is designed both for generality and portability across different Unix platforms, so you don’t have to rewrite your program each time you move to a new system. We’ve also made an effort to use reasonable, Scheme-like names for the multitudinous named constants involved, so when you are reading code, you’ll have less likelihood of getting lost in a bewildering maze of obfuscatory constants named ICRNL, INPCK, IUCLC, and ONOCR.

This section can only lay out the basic functionality of the terminal device interface. For further details, see the `termios(3)` man page on your system, or consult one of the standard Unix texts.

### 3.12.1 Portability across OS variants

Terminal-control software is inescapably complex, ugly, and low-level. Unix variants each provide their own way of controlling terminal devices, making it difficult to provide interfaces that are portable across different Unix systems. Scsh’s terminal support is based primarily upon the POSIX `termios` interface. Programs that can be written using only the POSIX interface are likely to be widely portable.

The bulk of the documentation that follows consists of several pages worth of tables defining different named constants that enable and disable different features of the terminal driver. Some of these flags are POSIX; others are taken from the two common branches of Unix development, SVR4 and 4.3+ Berkeley. Scsh guarantees that the non-POSIX constants will be bound identifiers.

- If your OS supports a particular non-POSIX flag, its named constant will be bound to the flag’s value.
- If your OS doesn’t support the flag, its named constant will be present, but bound to `#f`.

This means that if you want to use SVR4 or Berkeley features in a program, your program can portably test the values of the flags before using them—the flags can reliably be referenced without producing OS-dependent “unbound variable” errors.

Finally, note that although POSIX, SVR4, and Berkeley cover the lion’s share of the terminal-driver functionality, each operating system inevitably has non-standard extensions. While a particular scsh implementation may provide these extensions, they are not portable, and so are not documented here.

### 3.12.2 Miscellaneous procedures

`(tty? fd/port)`  $\longrightarrow$  *boolean* procedure  
Return true if the argument is a tty.

`(tty-file-name fd/port)`  $\longrightarrow$  *string* procedure  
The argument *fd/port* must be a file descriptor or port open on a tty. Return the file-name of the tty.

### 3.12.3 The tty-info record type

The primary data-structure that describes a terminal's mode is a `tty-info` record, defined as follows:

```
(define-record tty-info
  control-chars ; String: Magic input chars
  input-flags   ; Int: Input processing
  output-flags  ; Int: Output processing
  control-flags ; Int: Serial-line control
  local-flags   ; Int: Line-editting UI
  input-speed   ; Int: Code for input speed
  output-speed  ; Int: Code for output speed
  min           ; Int: Raw-mode input policy
  time)         ; Int: Raw-mode input policy
```

#### The control-characters string

The `control-chars` field is a character string; its characters may be indexed by integer values taken from table 3.4.

As discussed above, only the POSIX entries in table 3.4 are guaranteed to be legal, integer indices. A program can reliably test the OS to see if the non-POSIX characters are supported by checking the index constants. If the control-character function is supported by the terminal driver, then the corresponding index will be bound to an integer; if it is not supported, the index will be bound to `#f`.

To disable a given control-character function, set its corresponding entry in the `tty-info:control-chars` string to the special character `disable-tty-char` (and then use the `(set-tty-info fd/port info)` procedure to update the terminal's state).

## The flag fields

The `tty-info` record's `input-flags`, `output-flags`, `control-flags`, and `local-flags` fields are all bit sets represented as two's-complement integers. Their values are composed by or'ing together values taken from the named constants listed in tables 3.5 through 3.9.

As discussed above, only the POSIX entries listed in these tables are guaranteed to be legal, integer flag values. A program can reliably test the OS to see if the non-POSIX flags are supported by checking the named constants. If the feature is supported by the terminal driver, then the corresponding flag will be bound to an integer; if it is not supported, the flag will be bound to `#f`.

## The speed fields

The `input-speed` and `output-speed` fields determine the I/O rate of the terminal's line. The value of these fields is an integer giving the speed in bits-per-second. The following speeds are supported by POSIX:

0	134	600	4800
50	150	1200	9600
75	200	1800	19200
110	300	2400	38400

Your OS may accept others; it may also allow the special symbols `'exta` and `'extb`.

## The min and time fields

The integer `min` and `time` fields determine input blocking behaviour during non-canonical (raw) input; otherwise, they are ignored. See the `termios(3)` man page for further details.

Be warned that POSIX allows the base system call's representation of the `tty-info` record to share storage for the `min` field and the `ttychar/eof` element of the control-characters string, and for the `time` field and the `ttychar/eol` element of the control-characters string. Many implementations in fact do this.

To stay out of trouble, set the `min` and `time` fields only if you are putting the terminal into raw mode; set the `eof` and `eol` control-characters only if you are putting the terminal into canonical mode. It's ugly, but it's Unix.

### 3.12.4 Using tty-info records

(make-tty-info *if of cf lf ispeed ospeed min time*) → *tty-info-record* procedure  
 (copy-tty-info *tty-info-record*) → *tty-info-record* procedure

These procedures make it possible to create new `tty-info` records. The typical method for creating a new record is to copy one retrieved by a call to the `tty-info` procedure, then modify the copy as desired. Note that the `make-tty-info` procedure does not take a parameter to define the new record's control characters.<sup>6</sup> Instead, it simply returns a `tty-info` record whose control-character string has all elements initialised to ASCII nul. You may then install the special characters by assigning to the string. Similarly, the control-character string in the record produced by `copy-tty-info` does not share structure with the string in the record being copied, so you may mutate it freely.

(tty-info [*fd/port/fname*]) → *tty-info-record* procedure

The *fd/port/fname* parameter is an integer file descriptor or Scheme I/O port opened on a terminal device, or a file-name for a terminal device; it defaults to the current input port. This procedure returns a `tty-info` record describing the terminal's current mode.

(set-tty-info/now *fd/port/fname info*) → *no-value* procedure  
 (set-tty-info/drain *fd/port/fname info*) → *no-value* procedure  
 (set-tty-info/flush *fd/port/fname info*) → *no-value* procedure

The *fd/port/fname* parameter is an integer file descriptor or Scheme I/O port opened on a terminal device, or a file-name for a terminal device. The procedure chosen determines when and how the terminal's mode is altered:

Procedure	Meaning
set-tty-info/now	Make change immediately.
set-tty-info/drain	Drain output, then change.
set-tty-info/flush	Drain output, flush input, then change.

Oops: If I had defined these with the parameters in the reverse order, I could have made *fd/port/fname* optional. Too late now.

<sup>6</sup> Why? Because the length of the string varies from Unix to Unix. For example, the word-erase control character (typically control-w) is provided by most Unixes, but not part of the POSIX spec.

### 3.12.5 Other terminal-device procedures

(send-tty-break [*fd/port/fname* *duration*]) → *no-value* procedure

The *fd/port/fname* parameter is an integer file descriptor or Scheme I/O port opened on a terminal device, or a file-name for a terminal device; it defaults to the current output port. Send a break signal to the designated terminal. A break signal is a sequence of continuous zeros on the terminal's transmission line.

The *duration* argument determines the length of the break signal. A zero value (the default) causes a break of between 0.25 and 0.5 seconds to be sent; other values determine a period in a manner that will depend upon local community standards.

(drain-tty [*fd/port/fname*]) → *no-value* procedure

The *fd/port/fname* parameter is an integer file descriptor or Scheme I/O port opened on a terminal device, or a file-name for a terminal device; it defaults to the current output port.

This procedure waits until all the output written to the terminal device has been transmitted to the device. If *fd/port/fname* is an output port with buffered I/O enabled, then the port's buffered characters are flushed before waiting for the device to drain.

(flush-tty/input [*fd/port/fname*]) → *no-value* procedure

(flush-tty/output [*fd/port/fname*]) → *no-value* procedure

(flush-tty/both [*fd/port/fname*]) → *no-value* procedure

The *fd/port/fname* parameter is an integer file descriptor or Scheme I/O port opened on a terminal device, or a file-name for a terminal device; it defaults to the current input port (flush-tty/input and flush-tty/both), or output port (flush-tty/output).

These procedures discard the unread input chars or unwritten output chars in the tty's kernel buffers.

(start-tty-output [*fd/port/fname*]) → *no-value* procedure

(stop-tty-output [*fd/port/fname*]) → *no-value* procedure

(start-tty-input [*fd/port/fname*]) → *no-value* procedure

(stop-tty-input [*fd/port/fname*]) → *no-value* procedure

These procedures can be used to control a terminal's input and output flow. The *fd/port/fname* parameter is an integer file descriptor or Scheme I/O port opened on a terminal device, or a file-name for a terminal device; it defaults to the current input or output port.

The stop-tty-output and start-tty-output procedures suspend and resume output from a terminal device. The stop-tty-input and

start-tty-input procedures transmit the special STOP and START characters to the terminal with the intention of stopping and starting terminal input flow.

### 3.12.6 Control terminals, sessions, and terminal process groups

(open-control-tty *tty-name* [*flags*])  $\longrightarrow$  *port* procedure

This procedure opens terminal device *tty-name* as the process' control terminal (see the `termios` man page for more information on control terminals). The *tty-name* argument is a file-name such as `/dev/ttya`. The *flags* argument is a value suitable as the second argument to the `open-file` call; it defaults to `open/read+write`, causing the terminal to be opened for both input and output.

The port returned is an input port if the *flags* permit it, otherwise an output port. R5RS/Scheme 48/scsh do not have input/output ports, so it's one or the other. However, you can get both read and write ports open on a terminal by opening it `read/write`, taking the result input port, and duping it to an output port with `dup->outport`.

This procedure guarantees to make the opened terminal the process' control terminal only if the process does not have an assigned control terminal at the time of the call. If the `scsh` process already has a control terminal, the results are undefined.

To arrange for the process to have no control terminal prior to calling this procedure, use the `become-session-leader` procedure.

(become-session-leader)  $\longrightarrow$  *integer* procedure

This is the C `setsid()` call. POSIX job-control has a three-level hierarchy: session/process-group/process. Every session has an associated control terminal. This procedure places the current process into a brand new session, and disassociates the process from any previous control terminal. You may subsequently use `open-control-tty` to open a new control terminal.

It is an error to call this procedure if the current process is already a process-group leader. One way to guarantee this is not the case is only to call this procedure after forking.

(tty-process-group *fd/port/fname*)  $\longrightarrow$  *integer* procedure

(set-tty-process-group *fd/port/fname pgrp*)  $\longrightarrow$  *undefined* procedure

This pair of procedures gets and sets the process group of a given terminal.

(control-tty-file-name)  $\longrightarrow$  *string* procedure

Return the file-name of the process' control tty. On every version of Unix of which we are aware, this is just the string `"/dev/tty"`. However, this procedure uses the official Posix interface, so it is more portable than simply using a constant string.

### 3.12.7 Pseudo-terminals

Scsh implements an interface to Berkeley-style pseudo-terminals.

(fork-pty-session *thunk*)  $\longrightarrow$  [*process* *pty-in* *pty-out* *tty-name*] procedure

This procedure gives a convenient high-level interface to pseudo-terminals. It first allocates a pty/tty pair of devices, and then forks a child to execute procedure *thunk*. In the child process

- Stdio and the current I/O ports are bound to the terminal device.
- The child is placed in its own, new session (see `become-session-leader`).
- The terminal device becomes the new session's controlling terminal (see `open-control-tty`).
- The `(error-output-port)` is unbuffered.

The `fork-pty-session` procedure returns four values: the child's process object, two ports open on the controlling pty device, and the name of the child's corresponding terminal device.

(open-pty)  $\longrightarrow$  *pty-inport* *tty-name* procedure

This procedure finds a free pty/tty pair, and opens the pty device with read/write access. It returns a port on the pty, and the name of the corresponding terminal device.

The port returned is an input port—Scheme doesn't allow input/output ports. However, you can easily use `(dup->outport pty-inport)` to produce a matching output port. You may wish to turn off I/O buffering for this output port.

(pty-name->tty-name *pty-name*)  $\longrightarrow$  *tty-name* procedure

(tty-name->pty-name *tty-name*)  $\longrightarrow$  *pty-name* procedure

These two procedures map between corresponding terminal and pty controller names. For example,

(pty-name->tty-name `"/dev/ptyq3"`)  $\implies$  `"/dev/ttyq3"`  
(tty-name->pty-name `"/dev/ttyrc"`)  $\implies$  `"/dev/ptyrc"`

*Remark:* This is rather Berkeley-specific. SVR4 ptys are rare enough that I've no real idea if it generalises across the Unix gap. Experts are invited to advise. Users feel free to not worry—the predominance of current popular Unix systems use Berkeley ptys.

(make-pty-generator)  $\longrightarrow$  *procedure* procedure

This procedure returns a generator of candidate pty names. Each time the returned procedure is called, it produces a new candidate. Software that wishes to search through the set of available ptys can use a pty generator to iterate over them. After producing all the possible ptys, a generator returns #f every time it is called. Example:

```
(define pg (make-pty-generator))
(pg)  $\Rightarrow$  "/dev/ptyp0"
(pg)  $\Rightarrow$  "/dev/ptyp1"
      :
      :
(pg)  $\Rightarrow$  "/dev/ptyqe"
(pg)  $\Rightarrow$  "/dev/ptyqf"      (Last one)
(pg)  $\Rightarrow$  #f
(pg)  $\Rightarrow$  #f
      :
```



Scsh	C	Typical char
POSIX		
ttychar/delete-char	ERASE	del
ttychar/delete-line	KILL	^U
ttychar/eof	EOF	^D
ttychar/eol	EOL	
ttychar/interrupt	INTR	^C
ttychar/quit	QUIT	^\
ttychar/suspend	SUSP	^Z
ttychar/start	START	^Q
ttychar/stop	STOP	^S
SVR4 and 4.3+BSD		
ttychar/delayed-suspend	DSUSP	^Y
ttychar/delete-word	WERASE	^W
ttychar/discard	DISCARD	^O
ttychar/eol2	EOL2	
ttychar/literal-next	LNEXT	^V
ttychar/reprint	REPRINT	^R
4.3+BSD		
ttychar/status	STATUS	^T

Table 3.4: Indices into the `tty-info` record's *control-chars* string, and the character traditionally found at each index. Only the indices for the POSIX entries are guaranteed to be non-#f.

Scsh	C	Meaning
POSIX		
ttyin/check-parity	INPCK	Check parity.
ttyin/ignore-bad-parity-chars	IGNPAR	Ignore chars with parity errors.
ttyin/mark-parity-errors	PARMRK	Insert chars to mark parity errors.
ttyin/ignore-break	IGNBRK	Ignore breaks.
ttyin/interrupt-on-break	BRKINT	Signal on breaks.
ttyin/7bits	ISTRIP	Strip char to seven bits.
ttyin/cr->nl	ICRNL	Map carriage-return to newline.
ttyin/ignore-cr	IGNCR	Ignore carriage-returns.
ttyin/nl->cr	INLCR	Map newline to carriage-return.
ttyin/input-flow-ctl	IXOFF	Enable input flow control.
ttyin/output-flow-ctl	IXON	Enable output flow control.
SVR4 and 4.3+BSD		
ttyin/xon-any	IXANY	Any char restarts after stop.
ttyin/beep-on-overflow	IMAXBEL	Ring bell when queue full.
SVR4		
ttyin/lowercase	IUCLC	Map upper case to lower case.

Table 3.5: Input-flags. These are the named flags for the `tty-info` record's *input-flags* field. These flags generally control the processing of input chars. Only the POSIX entries are guaranteed to be non-#f.

Scsh	C	Meaning
POSIX		
ttyout/enable	OPOST	Enable output processing.
SVR4 and 4.3+BSD		
ttyout/nl->crnl	ONLCR	Map nl to cr-nl.
4.3+BSD		
ttyout/discard-eot	ONOEOT	Discard EOT chars.
ttyout/expand-tabs	OXTABS <sup>7</sup>	Expand tabs.
SVR4		
ttyout/cr->nl	OCRNL	Map cr to nl.
ttyout/nl-does-cr	ONLRET	Nl performs cr as well.
ttyout/no-col0-cr	ONOCR	No cr output in column 0.
ttyout/delay-w/fill-char	OFILL	Send fill char to delay.
ttyout/fill-w/del	OFDEL	Fill char is ASCII DEL.
ttyout/uppercase	OLCUC	Map lower to upper case.

Table 3.6: Output-flags. These are the named flags for the `tty-info` record's *output-flags* field. These flags generally control the processing of output chars. Only the POSIX entries are guaranteed to be non-#f.

	Value	Comment
Backspace delay	ttyout/bs-delay	Bit-field mask
	ttyout/bs-delay0	
	ttyout/bs-delay1	
Carriage-return delay	ttyout/cr-delay	Bit-field mask
	ttyout/cr-delay0	
	ttyout/cr-delay1	
	ttyout/cr-delay2	
	ttyout/cr-delay3	
Form-feed delay	ttyout/ff-delay	Bit-field mask
	ttyout/ff-delay0	
	ttyout/ff-delay1	
Horizontal-tab delay	ttyout/tab-delay	Bit-field mask
	ttyout/tab-delay0	
	ttyout/tab-delay1	
	ttyout/tab-delay2	
	ttyout/tab-delayx	
Newline delay	ttyout/nl-delay	Bit-field mask
	ttyout/nl-delay0	
	ttyout/nl-delay1	
Vertical tab delay	ttyout/vtab-delay	Bit-field mask
	ttyout/vtab-delay0	
	ttyout/vtab-delay1	
All	ttyout/all-delay	Total bit-field mask

Table 3.7: Delay constants. These are the named flags for the `tty-info` record's *output-flags* field. These flags control the output delays associated with printing special characters. They are non-POSIX, and have non-#f values only on SVR4 systems.

Scsh	C	Meaning
POSIX		
ttyc/char-size	CSIZE	Character size mask
ttyc/char-size5	CS5	5 bits
ttyc/char-size6	CS6	6 bits
ttyc/char-size7	CS7	7 bits
ttyc/char-size8	CS8	8 bits
ttyc/enable-parity	PARENB	Generate and detect parity.
ttyc/odd-parity	PARODD	Odd parity.
ttyc/enable-read	CREAD	Enable reception of chars.
ttyc/hup-on-close	HUPCL	Hang up on last close.
ttyc/no-modem-sync	LOCAL	Ignore modem lines.
ttyc/2-stop-bits	CSTOPB	Send two stop bits.
4.3+BSD		
ttyc/ignore-flags	CIGNORE	Ignore control flags.
ttyc/CTS-output-flow-ctl	CCTS_OFLOW	CTS flow control of output
ttyc/RTS-input-flow-ctl	CRTS_IFLOW	RTS flow control of input
ttyc/carrier-flow-ctl	MDMBUF	

Table 3.8: Control-flags. These are the named flags for the `tty-info` record's *control-flags* field. These flags generally control the details of the terminal's serial line. Only the POSIX entries are guaranteed to be non-#f.

Scsh	C	Meaning
POSIX		
ttyl/canonical	ICANON	Canonical input processing.
ttyl/echo	ECHO	Enable echoing.
ttyl/echo-delete-line	ECHOK	Echo newline after line kill.
ttyl/echo-nl	ECHONL	Echo newline even if echo is off.
ttyl/visual-delete	ECHOE	Visually erase chars.
ttyl/enable-signals	ISIG	Enable ^C, ^Z signalling.
ttyl/extended	IEXTEN	Enable extensions.
ttyl/no-flush-on-interrupt	NOFLSH	Don't flush after interrupt.
ttyl/ttou-signal	ITOSTOP	SIGTTOU on background output.
SVR4 and 4.3+BSD		
ttyl/echo-ctl	ECHOCTL	Echo control chars as “^X”.
ttyl/flush-output	FLUSHO	Output is being flushed.
ttyl/hardcopy-delete	ECHOPRT	Visual erase for hardcopy.
ttyl/reprint-unread-chars	PENDIN	Retype pending input.
ttyl/visual-delete-line	ECHOKE	Visually erase a line-kill.
4.3+BSD		
ttyl/alt-delete-word	ALTWERASE	Alternate word erase algorithm
ttyl/no-kernel-status	NOKERNINFO	No kernel status on ^T.
SVR4		
ttyl/case-map	XCASE	Canonical case presentation

Table 3.9: Local-flags. These are the named flags for the `tty-info` record's *local-flags* field. These flags generally control the details of the line-editing user interface. Only the POSIX entries are guaranteed to be non-#f.

## Chapter 4

# Networking

The Scheme Shell provides a BSD-style sockets interface. There is not an official standard for a network interface for scsh to adopt (this is the subject of the forthcoming Posix.8 standard). However, Berkeley sockets are a *de facto* standard, being found on most Unix workstations and PC operating systems.

It is fairly straightforward to add higher-level network protocols such as smtp, telnet, or http on top of the the basic socket-level support scsh provides. The Scheme Underground has also released a network library with many of these protocols as a companion to the current release of scsh. See this code for examples showing the use of the sockets interface.

### 4.1 High-level interface

For convenience, and to avoid some of the messy details of the socket interface, we provide a high level socket interface. These routines attempt to make it easy to write simple clients and servers without having to think of many of the details of initiating socket connections. We welcome suggested improvements to this interface, including better names, which right now are solely descriptions of the procedure's action. This might be fine for people who already understand sockets, but does not help the new networking programmer.

`(socket-connect protocol-family socket-type . args)`  $\longrightarrow$  `socket` procedure

`socket-connect` is intended for creating client applications. `protocol-family` is specified as either `protocol-family/internet` or `protocol-family/unix`. `socket-type` is specified as either `socket-type/stream` or `socket-type/datagram`. See `socket` for a more complete description of these terms.

The variable *args* list is meant to specify protocol family specific information. For Internet sockets, this consists of two arguments: a host name and a port number. For Unix sockets, this consists of a pathname.

`socket-connect` returns a socket which can be used for input and output from a remote server. See `socket` for a description of the *socket record*.

`(bind-listen-accept-loop protocol-family proc arg) → does-not-return procedure`

`bind-listen-accept-loop` is intended for creating server applications. *protocol-family* is specified as either the `protocol-family/internet` or `protocol-family/unix`. *proc* is a procedure of two arguments: a socket and a socket-address. *arg* specifies a port number for Internet sockets or a pathname for Unix sockets. See `socket` for a more complete description of these terms.

*proc* is called with a socket and a socket address each time there is a connection from a client application. The socket allows communications with the client. The socket address specifies the address of the remote client.

This procedure does not return, but loops indefinitely accepting connections from client programs.

`(bind-prepare-listen-accept-loop protocol-family prepare proc arg) → does-not-return procedure`

Same as `bind-listen-accept-loop` but runs the thunk *prepare* after binding the address and before entering the loop. The typical task of the *prepare* procedure is to change the user id from the superuser to some unprivileged id once the address has been bound.

## 4.2 Sockets

`(create-socket protocol-family type [protocol]) → socket procedure`  
`(create-socket-pair type) → [socket1 socket2] procedure`  
`(close-socket socket) → undefined procedure`

A socket is one end of a network connection. Three specific properties of sockets are specified at creation time: the protocol-family, type, and protocol.

The *protocol-family* specifies the protocol family to be used with the socket. This also determines the address family of socket addresses, which are described in more detail below. Scsh currently supports the Unix internal protocols and the Internet protocols using the following constants:

```
protocol-family/unspecified
protocol-family/unix
protocol-family/internet
```

The *type* specifies the style of communication. Examples that your operating system probably provides are stream and datagram sockets. Others may be available depending on your system. Typical values are:

```
socket-type/stream
socket-type/datagram
socket-type/raw
```

The *protocol* specifies a particular protocol to use within a protocol family and type. Usually only one choice exists, but it's probably safest to set this explicitly. See the protocol database routines for information on looking up protocol constants.

New sockets are typically created with `create-socket`. However, `create-socket-pair` can also be used to create a pair of connected sockets in the `protocol-family/unix` protocol-family. The value of a returned socket is a *socket record*, defined to have the following structure:

```
(define-record socket
  family                ; protocol family
  inport                ; input-port
  outport)              ; output-port
```

The family specifies the protocol family of the socket. The `inport` and `outport` fields are ports that can be used for input and output, respectively. For a stream socket, they are only usable after a connection has been established via `connect-socket` or `accept-connection`. For a datagram socket, `outport` can be immediately using `send-message`, and `inport` can be used after `bind` has created a local address.

`close-socket` provides a convenient way to close a socket's port. It is preferred to explicitly closing the `inport` and `outport` because using `close` on sockets is not currently portable across operating systems.

(`port->socket` *port* *protocol-family*)  $\longrightarrow$  *socket* procedure

This procedure turns *port* into a socket object. The port's underlying file descriptor must be a socket with protocol family *protocol-family*. `port->socket` applies `dup->inport` and `dup->outport` to *port* to create the ports of the socket object.

`port->socket` comes in handy for writing servers which run as children of `inetd`: after receiving a connection `inetd` creates a socket and passes it as standard input to its child.



### 4.3 Socket addresses

The format of a socket-address depends on the address family of the socket. Address-family-specific routines are provided to convert protocol-specific addresses to socket addresses. The value returned by these routines is a *socket-address record*, defined to have the following visible structure:

```
(define-record socket-address
  family)                                ; address family
```

The family is one of the following constants:

```
address-family/unspecified
address-family/unix
address-family/internet
```

(unix-address->socket-address *pathname*) → *socket-address* procedure  
unix-address->socket-address returns a *socket-address* based on the string *pathname*. There is a system dependent limit on the length of *pathname*.

(internet-address->socket-address *host-address* *service-port*) → *socket-address* procedure  
internet-address->socket-address returns a *socket-address* based on an integer *host-address* and an integer *service-port*. Besides being a 32-bit host address, an Internet host address can also be one of the following constants:

```
internet-address/any
internet-address/loopback
internet-address/broadcast
```

The use of internet-address/any is described below in bind-socket. internet-address/loopback is an address that always specifies the local machine. internet-address/broadcast is used for network broadcast communications.

For information on obtaining a host's address, see the host-info function.

(socket-address->unix-address *socket-address*) → *pathname* procedure  
(socket-address->internet-address *socket-address*) → [*host-address* *service-port*] procedure

The routines socket-address->internet-address and socket-address->unix-address return the address-family-specific addresses. Be aware that most implementations don't correctly return anything more than an empty string for addresses in the Unix address-family.

## 4.4 Socket primitives

The procedures in this section are presented in the order in which a typical program will use them. Consult a text on network systems programming for more information on sockets.<sup>1</sup> The last two tutorials are freely available as part of BSD. In the absence of these, your Unix manual pages for socket might be a good starting point for information.

`(connect-socket socket socket-address)`  $\longrightarrow$  *undefined* procedure

`connect-socket` sets up a connection from a *socket* to a remote *socket-address*. A connection has different meanings depending on the socket type. A stream socket must be connected before use. A datagram socket can be connected multiple times, but need not be connected at all if the remote address is specified with each `send-message`, described below. Also, datagram sockets may be disassociated from a remote address by connecting to a null remote address.

`(connect-socket-no-wait socket socket-address)`  $\longrightarrow$  *boolean* procedure

`(connect-socket-successful? socket)`  $\longrightarrow$  *boolean* procedure

Just like `connect-socket`, `connect-socket-no-wait` sets up a connection from a *socket* to a remote *socket-address*. Unlike `connect-socket`, `connect-socket-no-wait` does not block if it cannot establish the connection immediately. Instead it will return `#f` at once. In this case a subsequent `select` on the output port of the socket will report the output port as ready as soon as the operation system has established the connection or as soon as setting up the connection led to an error. Afterwards, the procedure `connect-socket-successful?` can be used to test whether the connection has been established successfully or not.

`(bind-socket socket socket-address)`  $\longrightarrow$  *undefined* procedure

`bind-socket` assigns a certain local *socket-address* to a *socket*. Binding a socket reserves the local address. To receive connections after binding the socket, use `listen-socket` for stream sockets and `receive-message` for datagram sockets.

Binding an Internet socket with a host address of `internet-address/any` indicates that the caller does not care to specify from which local network

---

<sup>1</sup> Some recommended ones are:

- “Unix Network Programming” by W. Richard Stevens
- “An Introductory 4.3BSD Interprocess Communication Tutorial.” (reprinted in UNIX Programmer’s Supplementary Documents Volume 1, PS1:7)
- “An Advanced 4.3BSD Interprocess Communication Tutorial.” (reprinted in UNIX Programmer’s Supplementary Documents Volume 1, PS1:8)

interface connections are received. Binding an Internet socket with a service port number of zero indicates that the caller has no preference as to the port number assigned.

Binding a socket in the Unix address family creates a socket special file in the file system that must be deleted before the address can be reused. See `delete-file`.

`(listen-socket socket backlog)`  $\longrightarrow$  *undefined* procedure

`listen-socket` allows a stream *socket* to start receiving connections, allowing a queue of up to *backlog* connection requests. Queued connections may be accepted by `accept-connection`.

`(accept-connection socket)`  $\longrightarrow$  [*new-socket socket-address*] procedure

`accept-connection` receives a connection on a *socket*, returning a new socket that can be used for this connection and the remote socket address associated with the connection.

`(socket-local-address socket)`  $\longrightarrow$  *socket-address* procedure

`(socket-remote-address socket)`  $\longrightarrow$  *socket-address* procedure

Sockets can be associated with a local address or a remote address or both. `socket-local-address` returns the local *socket-address* record associated with *socket*. `socket-remote-address` returns the remote *socket-address* record associated with *socket*.

`(shutdown-socket socket how-to)`  $\longrightarrow$  *undefined* procedure

`shutdown-socket` shuts down part of a full-duplex socket. The method of shutting down is specified by the *how-to* argument, one of:

`shutdown/receives`  
`shutdown/sends`  
`shutdown/sends+receives`

## 4.5 Performing input and output on sockets

`(receive-message socket length [flags])`  $\longrightarrow$  [*string-or-#f socket-address*] procedure

`(receive-message! socket string [start] [end] [flags])`  $\longrightarrow$  [*count-or-#f socket-address*] procedure

`(receive-message/partial socket length [flags])`  $\longrightarrow$  [*string-or-#f socket-address*] procedure

`(receive-message!/partial socket string [start] [end] [flags])`  $\longrightarrow$  [*count-or-#f socket-address*] procedure

`(send-message socket string [start] [end] [flags] [socket-address])`  $\longrightarrow$  *undefined* procedure

`(send-message/partial socket string [start] [end] [flags] [socket-address])`  $\longrightarrow$  *count* procedure

For most uses, standard input and output routines such as `read-string` and `write-string` should suffice. However, in some cases an extended

interface is required. The `receive-message` and `send-message` calls parallel the `read-string` and `write-string` calls with a similar naming scheme.

One additional feature of these routines is that `receive-message` returns the remote *socket-address* and `send-message` takes an optional remote *socket-address*. This allows a program to know the source of input from a datagram socket and to use a datagram socket for output without first connecting it.

All of these procedures take an optional *flags* field. This argument is an integer bit-mask, composed by or'ing together the following constants:

```
message/out-of-band
message/peek
message/dont-route
```

See `read-string` and `write-string` for a more detailed description of the arguments and return values.

## 4.6 Socket options

```
(socket-option socket level option) → value           procedure
(set-socket-option socket level option value) → undefined procedure
```

`socket-option` and `set-socket-option` allow the inspection and modification, respectively, of several options available on sockets. The *level* argument specifies what protocol level is to be examined or affected. A level of `level/socket` specifies the highest possible level that is available on all socket types. A specific protocol number can also be used as provided by `protocol-info`, described below.

There are several different classes of socket options. The first class consists of boolean options which can be either true or false. Examples of this option type are:

```
socket/debug
socket/accept-connect
socket/reuse-address
socket/keep-alive
socket/dont-route
socket/broadcast
socket/use-loop-back
socket/oob-inline
socket/use-privileged
socket/cant-signal
tcp/no-delay
```

Value options are another category of socket options. Options of this type are an integer value. Examples of this option type are:

```
socket/send-buffer
socket/receive-buffer
socket/send-low-water
socket/receive-low-water
socket/error
socket/type
ip/time-to-live
tcp/max-segment
```

A third option type specifies how long for data to linger after a socket has been closed. There is only one option of this type: `socket/linger`. It is set with either `#f` to disable it or an integer number of seconds to linger and returns a value of the same type upon inspection.

The fourth and final option type of this time is a timeout option. There are two examples of this option type: `socket/send-timeout` and `socket/receive-timeout`. These are set with a real number of microseconds resolution and returns a value of the same type upon inspection.

## 4.7 Database-information entries

```
(host-info name-or-socket-address)  → host-info      procedure
(network-info name-or-socket-address) → network-info or #f procedure
(service-info name-or-number [protocol-name]) → service-info or #f procedure
(protocol-info name-or-number) → protocol-info or #f procedure
```

`host-info` allows a program to look up a host entry based on either its string *name* or *socket-address*. The value returned by this routine is a *host-info record*, defined to have the following structure:

```
(define-record host-info
  name           ; Host name
  aliases        ; Alternative names
  addresses)     ; Host addresses
```

`host-info` could fail and raise an error for one of the following reasons:

```
herror/host-not-found
herror/try-again
herror/no-recovery
herror/no-data
herror/no-address
```

`network-info` allows a program to look up a network entry based on either its string *name* or *socket-address*. The value returned by this routine is a *network-info record*, defined to have the following structure:

```
(define-record network-info
  name                ; Network name
  aliases             ; Alternative names
  net)                ; Network number
```

`service-info` allows a program to look up a service entry based on either its string *name* or integer *port*. The value returned by this routine is a *service-info record*, defined to have the following structure:

```
(define-record service-info
  name                ; Service name
  aliases             ; Alternative names
  port                ; Port number
  protocol)           ; Protocol name
```

`protocol-info` allows a program to look up a protocol entry based on either its string *name* or integer *number*. The value returned by this routine is a *protocol-info record*, defined to have the following structure:

```
(define-record protocol-info
  name                ; Protocol name
  aliases             ; Alternative names
  number)             ; Protocol number)
```

`network-info`, `service-info` and `protocol-info` return `#f` if the specified entity was not found.

## Chapter 5

# Strings and characters

Strings are the basic communication medium for Unix processes, so a Unix programming environment must have reasonable facilities for manipulating them. Scsh provides a powerful set of procedures for processing strings and characters. Besides the the facilities described in this chapter, scsh also provides

- **Regular expressions (chapter 6)**

A complete regular-expression system.

- **Field parsing, delimited record I/O and the awk loop (chapter 8)**

These procedures let you read in chunks of text delimited by selected characters, and parse each record into fields based on regular expressions (for example, splitting a string at every occurrence of colon or white-space). The `awk` form allows you to loop over streams of these records in a convenient way.

- **The SRFI-13 string libraries**

This pair of libraries contains procedures that create, fold, iterate over, search, compare, assemble, cut, hash, case-map, and otherwise manipulate strings. They are provided by the `string-lib` and `string-lib-internals` packages, and are also available in the default `scsh` package.

More documentation on these procedures can be found at URLs

<http://srfi.schemers.org/srfi-13/srfi-13.html>

<http://srfi.schemers.org/srfi-13/srfi-13.txt>

- **The SRFI-14 character-set library**

This library provides a set-of-characters abstraction, which is frequently useful when searching, parsing, filtering or otherwise operating on

strings and character data. The SRFI is provided by the `char-set-lib` package; it's bindings are also available in the default `scsh` package.

More documentation on this library can be found at URLs  
<http://srfi.schemers.org/srfi-14/srfi-14.html>  
<http://srfi.schemers.org/srfi-14/srfi-14.txt>

## 5.1 Manipulating file names

These procedures do not access the file-system at all; they merely operate on file-name strings. Much of this structure is patterned after the `gnu emacs` design. Perhaps a more sophisticated system would be better, something like the pathname abstractions of `COMMON LISP` or `MIT Scheme`. However, being Unix-specific, we can be a little less general.

### 5.1.1 Terminology

These procedures carefully adhere to the `POSIX` standard for file-name resolution, which occasionally entails some slightly odd things. This section will describe these rules, and give some basic terminology.

A *file-name* is either the file-system root ("`/`"), or a series of slash-terminated directory components, followed by a file component. Root is the only file-name that may end in slash. Some examples:

File name	Dir components	File component
<code>src/des/main.c</code>	<code>("src" "des")</code>	<code>"main.c"</code>
<code>/src/des/main.c</code>	<code>("" "src" "des")</code>	<code>"main.c"</code>
<code>main.c</code>	<code>()</code>	<code>"main.c"</code>

Note that the relative filename `src/des/main.c` and the absolute filename `/src/des/main.c` are distinguished by the presence of the root component `"` in the absolute path.

Multiple embedded slashes within a path have the same meaning as a single slash. More than two leading slashes at the beginning of a path have the same meaning as a single leading slash—they indicate that the file-name is an absolute one, with the path leading from root. However, `POSIX` permits the OS to give special meaning to *two* leading slashes. For this reason, the routines in this section do not simplify two leading slashes to a single slash.

A file-name in *directory form* is either a file-name terminated by a slash, *e.g.*, `"/src/des/"`, or the empty string, `""`. The empty string corresponds to the current working directory, whose file-name is dot (`"."`). Working backwards from the append-a-slash rule, we extend the syntax of `POSIX` file-names to define the



empty string to be a file-name form of the root directory `"/`. (However, `"/` is also acceptable as a file-name form for root.) So the empty string has two interpretations: as a file-name form, it is the file-system root; as a directory form, it is the current working directory. Slash is also an ambiguous form: `/` is both a directory-form and a file-name form.

The directory form of a file-name is very rarely used. Almost all of the procedures in scsh name directories by giving their file-name form (without the trailing slash), not their directory form. So, you say `"/usr/include"`, and `."`, not `"/usr/include/"` and `""`. The sole exceptions are `file-name-as-directory` and `directory-as-file-name`, whose jobs are to convert back-and-forth between these forms, and `file-name-directory`, whose job it is to split out the directory portion of a file-name. However, most procedures that expect a directory argument will coerce a file-name in directory form to file-name form if it does not have a trailing slash. Bear in mind that the ambiguous case, empty string, will be interpreted in file-name form, *i.e.*, as root.

### 5.1.2 Procedures

(file-name-directory? <i>fname</i> )	→	<i>boolean</i>	procedure
(file-name-non-directory? <i>fname</i> )	→	<i>boolean</i>	procedure

These predicates return true if the string is in directory form, or file-name form (see the above discussion of these two forms). Note that they both return true on the ambiguous case of empty string, which is both a directory (current working directory), and a file name (the file-system root).

File name	...-directory?	...-non-directory?
"src/des"	#f	#t
"src/des/"	#t	#f
"/"	#t	#f
"."	#f	#t
""	#t	#t

(file-name-as-directory *fname*)  $\longrightarrow$  *string* procedure

Convert a file-name to directory form. Basically, add a trailing slash if needed:

```
(file-name-as-directory "src/des")  => "src/des/"
(file-name-as-directory "src/des/") => "src/des/"
```

., /, and "" are special:

```
(file-name-as-directory ".")      => ""
(file-name-as-directory "/" )    => "/"
(file-name-as-directory "")      => "/"
```

(directory-as-file-name *fname*) → *string* procedure

Convert a directory to a simple file-name. Basically, kill a trailing slash if one is present:

(directory-as-file-name "foo/bar/") ⇒ "foo/bar"

/ and "" are special:

(directory-as-file-name "/") ⇒ "/"

(directory-as-file-name "") ⇒ "." (*i.e.*, the cwd)

(file-name-absolute? *fname*) → *boolean* procedure

Does *fname* begin with a root or ~ component? (Recognising ~ as a home-directory specification is an extension of POSIX rules.)

(file-name-absolute? "/usr/shivers") ⇒ #t

(file-name-absolute? "src/des") ⇒ #f

(file-name-absolute? "~/src/des") ⇒ #t

Non-obvious case:

(file-name-absolute? "") ⇒ #t (*i.e.*, root)

(file-name-directory *fname*) → *string or false* procedure

Return the directory component of *fname* in directory form. If the file-name is already in directory form, return it as-is.

(file-name-directory "/usr/bdc") ⇒ "/usr/"

(file-name-directory "/usr/bdc/") ⇒ "/usr/bdc/"

(file-name-directory "bdc/.login") ⇒ "bdc/"

(file-name-directory "main.c") ⇒ ""

Root has no directory component:

(file-name-directory "/") ⇒ ""

(file-name-directory "") ⇒ ""

(file-name-nondirectory *fname*) → *string* procedure

Return non-directory component of *fname*.

```

(file-name-nondirectory "/usr/ian")    ⇒ "ian"
(file-name-nondirectory "/usr/ian/")   ⇒ ""
(file-name-nondirectory "ian/.login")  ⇒ ".login"
(file-name-nondirectory "main.c")      ⇒ "main.c"
(file-name-nondirectory "")            ⇒ ""
(file-name-nondirectory "/")           ⇒ "/"

```

(split-file-name *fname*) → *string list* procedure

Split a file-name into its components.

```

(split-file-name "src/des/main.c")
⇒      ("src" "des" "main.c")

(split-file-name "/src/des/main.c")
⇒      (" " "src" "des" "main.c")

(split-file-name "main.c")
⇒      ("main.c")

(split-file-name "/")
⇒      (" ")

```

(path-list->file-name *path-list* [*dir*]) → *string* procedure

Inverse of split-file-name.

```

(path-list->file-name '("src" "des" "main.c"))
⇒ "src/des/main.c"
(path-list->file-name '(" " "src" "des" "main.c"))
⇒ "/src/des/main.c"

```

Optional *dir* arg anchors relative path-lists:

```

(path-list->file-name '("src" "des" "main.c")
                      "/usr/shivers")
⇒ "/usr/shivers/src/des/main.c"

```

The optional *dir* argument is usefully (cwd).

(file-name-extension *fname*) → *string* procedure

Return the file-name's extension.

```

(file-name-extension "main.c")    ⇒ ".c"
(file-name-extension "main.c.old") ⇒ ".old"
(file-name-extension "/usr/shivers") ⇒ ""

```

Weird cases:

```
(file-name-extension "foo.")    => "."
(file-name-extension "foo..")   => "."
```

Dot files are not extensions:

```
(file-name-extension "/usr/shivers/.login") => ""
```

`(file-name-sans-extension fname)`  $\longrightarrow$  *string* procedure

Return everything but the extension.

```
(file-name-sans-extension "main.c")    => "main"
(file-name-sans-extension "main.c.old") => "main.c"
(file-name-sans-extension "/usr/shivers")
=> "/usr/shivers"
```

Weird cases:

```
(file-name-sans-extension "foo.")    => "foo"
(file-name-sans-extension "foo..")   => "foo."
```

Dot files are not extensions:

```
(file-name-sans-extension "/usr/shivers/.login")
=> "/usr/shivers/.login"
```

Note that appending the results of `file-name-extension` and `file-name-sans-extension` in all cases produces the original file-name.

`(parse-file-name fname)`  $\longrightarrow$  [*dir name extension*] procedure

Let *f* be `(file-name-nondirectory fname)`. This function returns the three values:

- `(file-name-directory fname)`
- `(file-name-sans-extension f)`
- `(file-name-extension f)`

The inverse of `parse-file-name`, in all cases, is `string-append`. The boundary case of `/` was chosen to preserve this inverse.

`(replace-extension fname ext)`  $\longrightarrow$  *string* procedure

This procedure replaces *fname*'s extension with *ext*. It is exactly equivalent to

```
(string-append (file-name-sans-extension fname) ext)
```

(simplify-file-name *fname*)  $\longrightarrow$  *string* procedure

Removes leading and internal occurrences of dot. A trailing dot is left alone, as the parent could be a symlink. Removes internal and trailing double-slashes. A leading double-slash is left alone, in accordance with POSIX. However, triple and more leading slashes are reduced to a single slash, in accordance with POSIX. Double-dots (parent directory) are left alone, in case they come after symlinks or appear in a `./../machine/...` “super-root” form (which POSIX permits).

(resolve-file-name *fname* [*dir*])  $\longrightarrow$  *string* procedure

- Do `~` expansion.
- If *dir* is given, convert a relative file-name to an absolute file-name, relative to directory *dir*.

(expand-file-name *fname* [*dir*])  $\longrightarrow$  *string* procedure

Resolve and simplify the file-name.

(absolute-file-name *fname* [*dir*])  $\longrightarrow$  *string* procedure

Convert file-name *fname* into an absolute file name, relative to directory *dir*, which defaults to the current working directory. The file name is simplified before being returned.

This procedure does not treat a leading tilde character specially.

(home-dir [*user*])  $\longrightarrow$  *string* procedure

home-dir returns *user*’s home directory. *User* defaults to the current user.

```
(home-dir)            $\implies$   "/user1/lecturer/shivers"  
(home-dir "ctkwan")  $\implies$   "/user0/research/ctkwan"
```

(home-file [*user*] *fname*)  $\longrightarrow$  *string* procedure

Returns file-name *fname* relative to *user*’s home directory; *user* defaults to the current user.

```
(home-file "man")       $\implies$   "/usr/shivers/man"  
(home-file "fcmlau" "man")  $\implies$   "/usr/fcmlau/man"
```

The general substitute-env-vars string procedure, defined in the previous section, is also frequently useful for expanding file-names.

## 5.2 Other string manipulation facilities

(substitute-env-vars *fname*)  $\longrightarrow$  *string* procedure

Replace occurrences of environment variables with their values. An environment variable is denoted by a dollar sign followed by alphanumeric chars and underscores, or is surrounded by braces.

```
(substitute-env-vars "$USER/.login")  
   $\implies$  "shivers/.login"  
(substitute-env-vars "${USER}_log")  $\implies$  "shivers_log"
```

## 5.3 ASCII encoding

(char->ascii *character*)  $\longrightarrow$  *integer* procedure

(ascii->char *integer*)  $\longrightarrow$  *character* procedure

These are identical to char->integer and integer->char except that they use the ASCII encoding.

## 5.4 Character predicates

(char-letter? *character*)  $\longrightarrow$  *boolean* procedure

(char-lower-case? *character*)  $\longrightarrow$  *boolean* procedure

(char-upper-case? *character*)  $\longrightarrow$  *boolean* procedure

(char-title-case? *character*)  $\longrightarrow$  *boolean* procedure

(char-digit? *character*)  $\longrightarrow$  *boolean* procedure

(char-letter+digit? *character*)  $\longrightarrow$  *boolean* procedure

(char-graphic? *character*)  $\longrightarrow$  *boolean* procedure

(char-printing? *character*)  $\longrightarrow$  *boolean* procedure

(char-whitespace? *character*)  $\longrightarrow$  *boolean* procedure

(char-blank? *character*)  $\longrightarrow$  *boolean* procedure

(char-iso-control? *character*)  $\longrightarrow$  *boolean* procedure

(char-punctuation? *character*)  $\longrightarrow$  *boolean* procedure

(char-hex-digit? *character*)  $\longrightarrow$  *boolean* procedure

(char-ascii? *character*)  $\longrightarrow$  *boolean* procedure

Each of these predicates tests for membership in one of the standard character sets provided by the SRFI-14 character-set library. Additionally, the following redundant bindings are provided for R5RS compatibility:

R5RS name	scsh definition
char-alphabetic?	char-letter+digit?
char-numeric?	char-digit?
char-alphanumeric?	char-letter+digit?

## 5.5 Deprecated character-set procedures

The SRFI-13 character-set library grew out of an earlier library developed for scsh. However, the SRFI standardisation process introduced incompatibilities with the original scsh bindings. The current version of scsh provides the library `obsolete-char-set-lib`, which contains the old bindings found in previous releases of scsh. The following table lists the members of this library, along with the equivalent SRFI-13 binding. This obsolete library is deprecated and *not* open by default in the standard scsh environment; new code should use the SRFI-13 bindings.

Old obsolete-char-set-lib	SRFI-13 char-set-lib
char-set-members	char-set->list
chars->char-set	list->char-set
ascii-range->char-set	ucs-range->char-set (not exact)
predicate->char-set	char-set-filter (not exact)
char-set-every?	char-set-every
char-set-any?	char-set-any
char-set-invert	char-set-complement
char-set-invert!	char-set-complement!
char-set:alphabetic	char-set:letter
char-set:numeric	char-set:digit
char-set:alphanumeric	char-set:letter+digit
char-set:control	char-set:iso-control

Note also that the `->char-set` procedure no longer handles a predicate argument.

## Chapter 6

# Pattern-matching strings with regular expressions

Scsh provides a rich facility for matching regular-expression patterns in strings. The system is composed of several pieces:

- An s-expression notation for writing down general regular expressions. In most systems, regexp patterns are encoded as string literals, such as "g(oo|ee)se". In scsh, they are written using s-expressions, such as (: "g" (| "oo" "ee") "se"), and are called *sre's*. The sre notation has several advantages over the traditional string-based notation. It's more expressive, can be commented, and can be indented to expose the structure of the form.
- An abstract data type (ADT) representation for regexp values. Traditional regular-expression systems compute regular expressions from run-time values using strings. This can be awkward. Scsh, instead, provides a separate data type for regexps, with a set of basic constructor and accessor functions; regular expressions can be dynamically computed and manipulated using these functions.
- Some tools that work on the regexp ADT: case-sensitive to case-insensitive regexp transform, a regexp simplifier, and so forth.
- Parsers and unparsers that can convert between external representations and the regexp ADT. The supported external representations are
  - Posix strings
  - S-expression notation (that is, sre's)



Being able to convert regexps to Posix strings allows implementations to implement regexp matching using standard Posix C-based engines.

- Macro support for the s-expression notation. The `rx` macro provides a new special form that allows you to embed regexps in the s-expression notation within a Scheme program. Evaluating the macro form produces a regexp ADT value which can be used by Scheme pattern-matching procedures and other regexp consumers.
- Pattern-matching and searching procedures. Spencer's Posix regexp engine is linked in to the runtime; the regexp code uses this engine to provide text matching.

The regexp language supported is a complete superset of Posix functionality, providing:

- sequencing and choice (`|`)
- repetition (`*`, `+`, `?`, `{m,n}`)
- character classes (e.g., `[aeiou]`) and wildcard (`.`)
- beginning/end of string anchors (`^`, `$`)
- case-sensitivity control
- submatch-marking

## 6.1 Summary SRE syntax

The following figures give a summary of the SRE syntax; the next section is a friendlier tutorial introduction.

<i>string</i>	Literal match—interpreted relative to the current case-sensitivity lexical context (default is case-sensitive)
<i>(string1 string2 ...)</i>	Set of chars, <i>e.g.</i> , ("abc" "XYZ"). Interpreted relative to the current case-sensitivity lexical context.
<i>(* sre ...)</i>	0 or more matches
<i>(+ sre ...)</i>	1 or more matches
<i>(? sre ...)</i>	0 or 1 matches
<i>(= n sre ...)</i>	<i>n</i> matches
<i>(&gt;= n sre ...)</i>	<i>n</i> or more matches
<i>(** n m sre ...)</i>	<i>n</i> to <i>m</i> matches
	<i>N</i> and <i>m</i> are Scheme expressions producing non-negative integers. <i>M</i> may also be #f, meaning “infinity.”
<i>(  sre ...)</i>	Choice (or is R5RS symbol;
<i>(or sre ...)</i>	is not specified by R5RS.)
<i>(: sre ...)</i>	Sequence (seq is legal
<i>(seq sre ...)</i>	Common Lisp symbol)
<i>(submatch sre ...)</i>	Numbered submatch
<i>(dsm pre post sre ...)</i>	Deleted submatches
	<i>Pre</i> and <i>post</i> are numerals.
<i>(uncase sre ...)</i>	Case-folded match
<i>(w/case sre ...)</i>	Introduce a lexical case-sensitivity
<i>(w/nocase sre ...)</i>	context.
<i>,@exp</i>	Dynamically computed regexp
<i>,exp</i>	Same as <i>,@exp</i> , but no submatch info
	<i>Exp</i> must produce a character, string, char-set, or regexp.
<i>bos eos</i>	Beginning/end of string
<i>bol eol</i>	Beginning/end of line

Figure 6.1: SRE syntax summary (part 1)

<code>(posix-string string)</code>	Escape for Posix string notation
<code>char</code>	Singleton char set
<code>class-name</code>	alphanumeric, whitespace, etc.
These two forms are interpreted subject to the lexical case-sensitivity context.	
<code>(~ cset-sre ...)</code>	Complement-of-union ( <code>[^...]</code> )
<code>(- cset-sre ...)</code>	Difference
<code>(&amp; cset-sre ...)</code>	Intersection
<code>(/ range-spec ...)</code>	Character range—interpreted subject to the lexical case-sensitivity context

Figure 6.2: SRE syntax summary (part 2)

<code>class-name</code>	<code>::=</code>	any
		nonl
		lower-case   lower
		upper-case   upper
		alphabetic   alpha
		numeric   digit   num
		alphanumeric   alnum
		punctuation   punct
		graphic   graph
		whitespace   space   white
		printing   print
		control   cntrl
		hex-digit   xdigit   hex
		ascii
<code>range-spec</code>	<code>::=</code>	<code>string</code>   <code>char</code>
The chars are taken in pairs to form inclusive ranges.		

Figure 6.3: SRE character-class names and range specs.

<code>&lt;cset-sre&gt; ::= (~ &lt;cset-sre&gt; ...)</code>	Set complement-of-union
<code>  (- &lt;cset-sre&gt; ...)</code>	Set difference
<code>  (&amp; &lt;cset-sre&gt; ...)</code>	Intersection
<code>  (  &lt;cset-sre&gt; ...)</code>	Set union
<code>  (/ &lt;range-spec&gt; ...)</code>	Range
<code>  (&lt;string&gt;)</code>	Constant set
<code>  &lt;char&gt;</code>	Singleton constant set
<code>  &lt;string&gt;</code>	For 1-char string "c"
<code>  &lt;class-name&gt;</code>	Constant set
<code>  ,&lt;exp&gt;</code>	<code>&lt;exp&gt;</code> evals to a char-set,
<code>  ,@&lt;exp&gt;</code>	char, single-char string,
	or re-char-set regexp.
<code>  (uncase &lt;cset-sre&gt;)</code>	Case-folding
<code>  (w/case &lt;cset-sre&gt;)</code>	
<code>  (w/nocase &lt;cset-sre&gt;)</code>	

Figure 6.4: applied to SRE's that specify character sets. These are the "type-checking" rules for character-set SRE's.

## 6.2 Examples

```
(- alpha ("aeiouAEIOU"))          ; Various forms of
(- alpha ("aeiou") ("AEIOU"))      ; non-vowel letter
(w/nocase (- alpha ("aeiou")))
(- (/ "azAZ") ("aeiouAEIOU"))
(w/nocase (- (/ "az") ("aeiou")))

;;; Upper-case letter, lower-case vowel, or digit
(| upper ("aeiou") digit)
(| (/ "AZ09") ("aeiou"))

;;; Not an SRE, but Scheme code containing some embedded SREs.
(let* ((ws (rx (+ whitespace)))      ; Seq of whitespace
      (date (rx (: (| "Jan" "Feb" "Mar" ...) ; A month/day date.
                    ,ws
                    (| ("123456789")          ; 1-9
                      (: ("12") digit)        ; 10-29
                      "30" "31")))))          ; 30-31

      ;; Now we can use DATE several times:
      (rx ... ,date ... (* ... ,date ...)
        ... .... ,date))

;;; More Scheme code
(define (csl re)                      ; A comma-separated list of RE's is
  (rx (| ""                          ; either zero of them (empty string), or
      (: ,re                          ; one RE, followed by
        (* " , " ,re))))             ; Zero or more comma-space-RE matches.

  (csl (rx (| "John" "Paul" "George" "Ringo"))))
```

## 6.3 A short tutorial

S-expression regexps are called "SRE"s. Keep in mind that they are *not* Scheme expressions; they are another, separate notation that is expressed using the underlying framework of s-expression list structure: lists, symbols, *etc.* SRE's can be *embedded* inside of Scheme expressions using special forms that extend Scheme's syntax (such as the `rx` macro); there are places in the SRE grammar where one may place a Scheme expression. In these ways, SRE's and Scheme expressions can be intertwined. But this isn't fundamental; SRE's may be used in a completely Scheme-independent context. By simply restricting the nota-

tion to eliminate two special Scheme-embedding forms, they can be a completely independent notation.

**Constant strings** The simplest SRE is a string, denoting a constant regexp. For example, the SRE

```
"Spot"
```

matches only the string <<capital-S, little-p, little-o, little-t>>. There is no interpretation of the characters in the string at all—the SRE

```
".*["
```

matches the string <<period, asterisk, open-bracket>>.

**Simple character sets** To specify a set of characters, write a list whose single element is a string containing the set’s elements. So the SRE

```
("aeiou")
```

only matches a vowel. One way to think of this, notationally, is that the set brackets are "(" and ")".

**Wild card** Another simple SRE is the symbol `any`, which matches any single character—including newline and ASCII `nul`.

**Sequences** We can form sequences of SRE’s with the SRE `(: sre ...)`. So the SRE

```
(: "x" any "z")
```

matches any three-character string starting with “x” and ending with “z”. As we’ll see shortly, many SRE forms have bodies that are implicit sequences of other SRE’s, analogous to the manner in which the body of a Scheme `lambda` or `let` expression is an implicit `begin` sequence. The regexp `(seq sre ...)` is completely equivalent to `(: sre ...)`; it’s included in order to have a syntax that doesn’t require `:` to be a legal symbol<sup>1</sup>

---

<sup>1</sup>That is, for use within s-expression syntax frameworks that, unlike R5RS, don’t allow for `:` as a legal symbol. A Common Lisp embedding of SREs, for example, would need to use `seq` instead of `:`.

## 6.4 Choices

The SRE (`| sre ...`) is a regexp that matches anything any of the *sre* regexps match. So the regular expression

```
(| "sasha" "Pete")
```

matches either the string “sasha” or the string “Pete”. The regexp

```
(| ("aeiou") ("0123456789"))
```

is the same as

```
("aeiou0123456789")
```

The regexp (`or sre ...`) is completely equivalent to (`| sre ...`); it’s included in order to have a syntax that doesn’t require `|` to be a legal symbol.

**Repetition** There are several SRE forms that match multiple occurrences of a regular expression. For example, the SRE (`* sre ...`) matches zero or more occurrences of the sequence (`: sre ...`). Here is the complete list of SRE repetition forms:

SRE	means	at least	no more than
( <code>* sre ...</code> )	zero-or-more	0	infinity
( <code>+ sre ...</code> )	one-or-more	1	infinity
( <code>? sre ...</code> )	zero-or-one	0	1
( <code>= from sre ...</code> )	exactly-n	<i>from</i>	<i>from</i>
( <code>&gt;= from sre ...</code> )	n-or-more	<i>from</i>	infinity
( <code>** from to sre ...</code> )	n-to-m	<i>from</i>	<i>to</i>

A *from* field is a Scheme expression that produces an integer. A *to* field is a Scheme expression that produces either an integer, or false, meaning infinity.

While it is illegal for the *from* or *to* fields to be negative, it *is* allowed for *from* to be greater than *to* in a `**` form—this simply produces a regexp that will never match anything.

As an example, we can describe the names of `car/cdr` access functions (“`car`”, “`cdr`”, “`cadr`”, “`cdar`”, “`caar`”, “`cddr`”, “`caadr`”, *etc.*) with either of the SREs

```
(: "c" (+ (| "a" "d")) "r")  
(: "c" (+ ("ad")) "r")
```

We can limit the a/d chains to 4 characters or less with the SRE

```
(: "c" (** 1 4 ("ad")) "r")
```

Some boundary cases:

```
(** 5 2 "foo")      ; Will never match
(** 0 0 "foo")      ; Matches the empty string
```

**Character classes** There is a special set of SRE's that form “character classes”—basically, a regexp that matches one character from some specified set of characters. There are operators to take the intersection, union, complement, and difference of character classes to produce a new character class. (Except for union, these capabilities are not provided for general regexps as they are computationally intractable in the general case.)

A single character is the simplest character class: `#\x` is a character class that matches only the character “x”. A string that has only one letter is also a character class: `"x"` is the same SRE as `#\x`.

The character-set notation (*string*) we've seen is a primitive character class, as is the wildcard `any`. When arguments to the choice operator, `|`, are all character classes, then the choice form is itself a character-class. So these SREs are all character-classes:

```
("aeiou")
(| #\a #\e #\i #\o #\u)
(| ("aeiou") ("1234567890"))
```

However, these SRE's are *not* character-classes:

```
"aeiou"
(| "foo" #\x)
```

The (`~ cset-sre ...`) char class matches one character not in the specified classes:

```
(~ ("0248") ("1359"))
```

matches any character that is not a digit.

More compactly, we can use the `/` operator to specify character sets by giving the endpoints of contiguous ranges, where the endpoints are specified by a sequence of strings and characters. For example, any of these char classes

```
(/ #\A #\Z #\a #\z #\0 #\9)
(/ "AZ" #\a #\z "09")
(/ "AZ" #\a "z09")
(/ "AZaz09")
```



matches a letter or a digit. The range endpoints are taken in pairs to form inclusive ranges of characters. Note that the exact set of characters included in a range is dependent on the underlying implementation's character type, so ranges may not be portable across different implementations.

There is a wide selection of predefined, named character classes that may be used. One such SRE is the wildcard `any`. `nonl` is a character class matching anything but newline; it is equivalent to

```
(~ #\newline)
```

and is useful as a wildcard in line-oriented matching.

There are also predefined named char classes for the standard Posix and Gnu character classes:

scsh name	Posix/ctype	Alternate name	Comment
lower-case	lower		
upper-case	upper		
alphabetic	alpha		
numeric	digit	num	
alphanumeric	alnum	alphanum	
punctuation	punct		
graphic	graph		
blank	(Gnu extension)		
whitespace	space	white	"space" is deprecated.
printing	print		
control	cntrl		
hex-digit	xdigit	hex	
ascii	(Gnu extension)		

See the `scsh` character-set documentation or the Posix `isalpha(3)` man page for the exact definitions of these sets.

You can use either the long `scsh` name or the shorter Posix and alternate names to refer to these char classes. The standard Posix name "space" is provided, but deprecated, since it is ambiguous. It means "whitespace," the set of whitespace characters, not the singleton set of the `#\space` character. If you want a short name for the set of whitespace characters, use the char-class name "white" instead.

Char classes may be intersected with the operator `(& cset-sre ...)`, and set-difference can be performed with `(- cset-sre ...)`. These operators are particularly useful when you want to specify a set by negation *with respect to a limited universe*. For example, the set of all non-vowel letters is

```
(- alpha ("aeiou") ("AEIOU"))
```

whereas writing a simple complement

```
(~ ("aeiouAEIOU"))
```

gives a char class that will match any non-vowel—including punctuation, digits, white space, control characters, and ASCII nul.

We can *compute* a char class by writing the SRE

```
,cset-exp
```

where *cset-exp* is a Scheme expression producing a value that can be coerced to a character set: a character set, character, one-character string, or char-class regexp value. This regexp matches one character from the set.

The char-class SRE `,@cset-exp` is entirely equivalent to `,cset-exp` when *cset-exp* produces a character set (but see below for the more general non-char-class context, where there *is* a distinction between `,exp` and `,@exp`).

As an example of character-class SREs, an SRE that matches a lower-case vowel, upper-case letter, or digit is

```
(| ("aeiou") ("/AZ09"))
```

or, equivalently

```
(| ("aeiou") upper-case numeric)
```

Boundary cases: the empty-complement char class

```
(~)
```

matches any character; it is equivalent to any. The empty-union char class

```
(|)
```

never matches at all. This is rarely useful for human-written regexps, but may be of occasional utility in machine-generated regexps, perhaps produced by macros.

The rules for determining if an SRE is a simple, char-class SRE or a more complex SRE form a little “type system” for SRE’s. See the summary section preceding this one for a complete listing of these rules.

{Note There is no way to include the ASCII NUL character in a character set or search for it in any other way using regular expression. This is because the POSIX regexp facility is based on the C language which uses ASCII NUL to terminate strings.}

**Case sensitivity** There are three forms that control case sensitivity:

```
(uncase  sre ...)  
(w/case  sre ...)  
(w/nocase sre ...)
```

`uncase` is a regexp operator producing a regexp that matches any case permutation of any string that matches `(: sre ...)`. For example, the regexp

```
(uncase "foo")
```

matches the strings `"foo"`, `"foO"`, `"fOo"`, `"fOO"`, `"Foo"`, ...

Expressions in SRE notation are interpreted in a lexical case-sensitivity context. The forms `w/case` and `w/nocase` are the scoping operators for this context, which controls how constant strings and char-class forms are interpreted in their bodies. So, for example, the regexp

```
(w/nocase "abc"  
  (* "F00" (w/case "Bar"))  
  ("aeiou"))
```

defines a case-insensitive match for all of its elements except for the subelement `"Bar"`, which must match exactly capital-B, little-a, little-r. The default, the outermost, top-level context is case sensitive.

The lexical case-sensitivity context affects the interpretation of

- constant strings, such as `"foo"`,
- chars, such as `#\x`,
- char sets, such as `("abc")`, and
- ranges, such as `(/"az")` that appear within that context. It does not affect dynamically computed regexps—ones that are introduced by `,exp` and `,@exp` forms. It does not affect named char-classes—presumably, if you wrote `lower`, you didn't mean `alpha`.

`uncase` is *not* the same as `w/nocase`. To point up one distinction, consider the two regexps

```
(uncase (~ "a"))  
(w/nocase (~ "a"))
```

The regexp (`~ "a"`) matches any character except "a," which means it *does* match "A." Now, (`uncase re`) matches any case-permutation of a string that *re* matches. (`~ "a"`) matches "A," so (`uncase (~ "a")`) matches "A" and "a"—and, for that matter, every other character. So (`uncase (~ "a")`) is equivalent to any.

In contrast, (`w/nocase (~ "a")`) establishes a case-insensitive lexical context in which the "a" is interpreted, making the SRE equivalent to (`~ ("aA")`).

**Dynamic regexps** SRE notation allows you to compute parts of a regular expressions at run time. The SRE

`,exp`

is a regexp whose body *exp* is a Scheme expression producing a string, character, char-set, or regexp as its value. Strings and characters are converted into constant regexps; char-sets are converted into char-class regexps; and regexp values are substituted in place. So we can write regexps like this

```
(: "feeding the "
, (if (> n 1) "geese" "goose"))
```

This is how you can drop computed strings, such as someone's name, or the decimal numeral for a computed number, into a complex regexp.

If we have a large, complex regular expression that is used multiple times in some other, containing regular expression, we can name it, using the binding forms of the embedding language (*e.g.*, Scheme), and refer to it by name in the containing expression. For example, consider the Scheme expression

```
(let* ((ws (rx (+ whitespace))) ; Seq of whitespace
      ; Something like "Mar 14"
      (date (rx (: (| "Jan" "Feb" "Mar" ...)
                    ,ws
                    (| ("123456789") ; 1-9
                        (: ("12") digit) ; 10-29
                        "30" ; 30
                        "31"))))) ; 31
      ; Now we can use DATE several times:
      (rx ... ,date ... (* ... ,date ...
                        ... ,date ...)))
```

where the (`rx sre ...`) macro is the Scheme special form that produces a Scheme regexp value given a body in SRE notation.

As we saw in the char-class section, if a dynamic regexp is used in a char-class context (*e.g.*, as an argument to a `~` operation), the expression must be

coercable not merely to a general regexp, but to a character sre—so it must be either a singleton string, a character, a scsh char set, or a char-class regexp.

We can also define and use functions on regexps in the host language. For example, consider the following Scheme expressions, containing embedded SRE's (inside the `rx` macro expressions) which in turn contain embedded Scheme expressions computing dynamic regexps:

```
(define (csl re)
  ;; A comma-separated list of RE's is either
  (rx (| "" ; zero of them (empty string),
        (: ,re ; or RE followed by
          (* " , " ,re))))); zero or more comma-space-RE matches.

  (rx ... ,date ...
    ,(csl (rx (| "John" "Paul" "George" "Ringo"))))
    ...
    ,(csl date)
    ...)
```

We leave the extension of `csl` to allow for an optional “and” between the last two matches as an exercise for the interested reader (*e.g.*, to match “John, Paul, George and Ringo”).

Note, in passing, one of the nice features of SRE notation: they can be commented, and indented in a fashion to show the lexical extent of the subexpressions.

When we embed a computed regexp inside another regular expression with the `,exp` form, we must specify how to account for the submatches that may be in the computed part. For example, suppose we have the regexp

```
(rx (submatch (* "foo"))
    (submatch (? "bar"))
    ,(f x)
    (submatch "baz"))
```

It's clear that the submatch for the `(* "foo")` part of the regexp is submatch #1, and the `(? "bar")` part is submatch #2. But what number submatch is the "baz" submatch? It's not clear. Suppose the Scheme expression `(f x)` produces a regular expression that itself has 3 subforms. Are these counted (making the "baz" submatch #6), or not counted (making the "bar" submatch #3)?

SRE notation provides for both possibilities. The SRE

```
,exp
```

does *not* contribute its submatches to its containing regexp; it has zero submatches. So one can reliably assign submatch indices to forms appearing after a *,exp* form in a regexp.

On the other hand, the SRE

*,@exp*

“splices” its resulting regexp into place, *exposing* its submatches to the containing regexp. This is useful if the computed regexp is defined to produce a certain number of submatches—if that is part of *exp*’s “contract.”

**String and line units** The regexps *bos* and *eos* match the empty string at the beginning and end of the string, respectively.

The regexps *bol* and *eor* match the empty string at the beginning and end of a line, respectively. A line begins at the beginning of the string, and just after every newline character. A line ends at the end of the string, and just before every newline character. The char class *nonl* matches any character except newline, and is useful in conjunction with line-based pattern matching.

{Note *bol* and *eor* are not supported by *scsh*’s current regexp search engine, which is Spencer’s Posix matcher. This is the only element of the notation that is not supported by the current *scsh* reference implementation.}

**Posix string notation** The SRE (*posix-string string*), where *string* is a string literal (*not* a general Scheme expression), allows one to use Posix string notation for a regexp. It’s intended as backwards compatibility and is deprecated. For example, (*posix-string* "[aeiou]+|x\*|y{3,5}") matches a string of vowels, a possibly empty string of x’s, or three to five y’s.

Note that parentheses are used ambiguously in Posix notation—both for grouping and submatch marking. The (*posix-string string*) form makes the conservative assumption: all parentheses introduce submatches.

**Deleted submatches** Deleted submatches, or “DSM’s,” are a subtle feature that are never required in expressions written by humans. They can be introduced by the simplifier when reducing regular expressions to simpler equivalents, and are included in the syntax to give it expressibility spanning the full regexp ADT. They may appear when unparsing simplified regular expressions that have been run through the simplifier; otherwise you are not likely to see them. Feel free to skip this section.

The regexp simplifier can sometimes eliminate entire sub-expressions from a regexp. For example, the regexp

```
(: "foo" (** 0 0 "apple") "bar")
```

can be simplified to

```
"foobar"
```

since `(** 0 0 "apple")` will always match the empty string. The regexp

```
(| "foo"  
  (: "Richard" (|) "Nixon")  
  "bar")
```

can be simplified to

```
(| "foo" "bar")
```

The empty choice `(|)` can't match anything, so the whole

```
(: "Richard" (|) "Nixon")
```

sequence can't match, and we can remove it from the choice.

However, if deleting part of a regular expression removes a submatch form, any following submatch forms will have their numbering changed, which would be an error. For example, if we simplify

```
(: (** 0 0 (submatch "apple"))  
  (submatch "bar"))
```

to

```
(submatch "bar")
```

then the "bar" submatch changes from submatch #2 to submatch #1—so this is not a legal simplification.

When the simplifier deletes a sub-regexp that contains submatches, it introduces a special regexp form to account for the missing, deleted submatches, thus keeping the submatch accounting correct.

```
(dsm pre post sre ...)
```

is a regexp that matches the sequence `(: sre ...)`. *pre* and *post* are integer constants. The DSM form introduces *pre* deleted submatches before the body, and *post* deleted submatches after the body. If the body `(: sre ...)` itself has *body-sm* submatches, then the total number of submatches for the DSM form is

$$pre + body-sm + post.$$

These extra, deleted submatches are never assigned string indices in any match values produced when matching the regexp against a string.

As examples,

```
(| (: (submatch "Richard") (|) "Nixon")
  (submatch "bar"))
```

can be simplified to

```
(dsm 1 0 (submatch "bar"))
```

The regexp

```
(: (** 0 0 (submatch "apple")
  (submatch "bar"))
```

can be simplified to

```
(dsm 1 0 (submatch "bar"))
```

#### 6.4.1 Embedding regexps within Scheme programs

SRE's can be placed in a Scheme program using the `(rx sre ...)` Scheme form, which evaluates to a Scheme regexp value.

##### Static and dynamic regexps

We separate SRE expressions into two classes: static and dynamic expressions. A *static* expression is one that has no run-time dependencies; it is a complete, self-contained description of a regular set. A *dynamic* expression is one that requires run-time computation to determine the particular regular set being described. There are two places where one can embed run-time computations in an SRE:

- The *from* or *to* repetition counts of `**`, `=`, and `>=` forms;
- `,exp` and `,@exp` forms.

A static SRE is one that does not contain any `,exp` or `,@exp` forms, and whose `**`, `=`, and `>=` forms all contain constant repetition counts.

Scsh's `rx` macro is able, at macro-expansion time, to completely parse, simplify and translate any static SRE into the equivalent Posix string which is used to drive the underlying C-based matching engine; there is no run-time overhead. Dynamic SRE's are partially simplified and then expanded into Scheme code that constructs the regexp at run-time.



## 6.5 Regexp functions

### 6.5.1 Obsolete, deprecated procedures

These two procedures are survivors from the previous, now-obsolete `scsh` regexp interface. Old code must open the `re-old-funs` package to access them. They should not be used in new code.

`(string-match posix-re-string string [start])`  $\longrightarrow$  *match or false* procedure  
`(make-regexp posix-re-string)`  $\longrightarrow$  *regexp* procedure

These are old functions included for backwards compatibility with previous releases. They are deprecated and will go away at some point in the future.

Note that the new release has no “regexp compiling” procedure at all—regexp values are compiled for the matching engine on-demand, and the necessary data structures are cached inside the ADT values.

### 6.5.2 Standard procedures and syntax

`(rx sre ...)`  $\longrightarrow$  *regexp* Syntax

This allows you to describe a regexp value with SRE notation.

`(regexp? x)`  $\longrightarrow$  *boolean* procedure

Returns true if the value is a regular expression.

`(regexp-search re string [start flags])`  $\longrightarrow$  *match-data or false* procedure  
`(regexp-search? re string [start flags])`  $\longrightarrow$  *boolean* procedure

Search *string* starting at position *start*, looking for a match for regexp *re*. If a match is found, return a match structure describing the match, otherwise `#f`. *Start* defaults to 0.

*Flags* is the bitwise-or of `regexp/bos-not-bol` and `regexp/eos-not-eol`. `regexp/bos-not-bol` means the beginning of the string isn’t a line-begin. `regexp/eos-not-eol` is analogous. {Note They’re currently ignored because begining/end-of-line anchors aren’t supported by the current implementation.}

Use `regexp-search?` when you don’t need submatch information, as it has the potential to be *significantly* faster on submatch-containing regexps.

There is no longer a separate regexp “compilation” function; regexp values are compiled for the C engine on demand, and the resulting C structures are cached in the regexp structure after the first use.

<code>(match:start m [i])</code>	$\longrightarrow$	<i>integer or false</i>	procedure
<code>(match:end m [i])</code>	$\longrightarrow$	<i>integer or false</i>	procedure
<code>(match:substring m [i])</code>	$\longrightarrow$	<i>string or false</i>	procedure

`match:start` returns the start position of the submatch denoted by *match-number*. The whole regexp is 0; positive integers index submatches in the regexp, counting left-to-right. *Match-number* defaults to 0.

If the regular expression matches as a whole, but a particular sub-expression does not match, then `match:start` returns `#f`.

`match:end` is analogous to `match:start`, returning the end position of the indexed submatch.

`match:substring` returns the substring matched regexp's submatch. If there was no match for the indexed submatch, it returns `false`.

`(regexp-substitute port-or-false match . items)`  $\longrightarrow$  *object*      procedure

This procedure can be used to perform string substitutions based on regular-expression matches. The results of the substitution can be either output to a port or returned as a string.

The *match* argument is a regular-expression match structure that controls the substitution. If *port* is an output port, the *items* are written out to the port:

- If an item is a string, it is copied directly to the port.
- If an item is an integer, the corresponding submatch from *match* is written to the port.
- If an item is `'pre`, the prefix of the matched string (the text preceding the match) is written to the port.
- If an item is `'post`, the suffix of the matched string is written.

If *port* is `#f`, nothing is written, and a string is constructed and returned instead.

`(regexp-substitute/global port-or-false re str . items)`  $\longrightarrow$  *object*      procedure

This procedure is similar to `regexp-substitute`, but can be used to perform repeated match/substitute operations over a string. It has the following differences with `regexp-substitute`:

- It takes a regular expression and string to be matched as parameters, instead of a completed match structure.
- If the regular expression doesn't match the string, this procedure is the identity transform—it returns or outputs the string.

- If an item is 'post, the procedure recurses on the suffix string (the text from *string* following the match). Including a 'post in the list of items is how one gets multiple match/substitution operations.
- If an item is a procedure, it is applied to the match structure for a given match. The procedure returns a string to be used in the result.

The *regexp* parameter can be either a compiled regular expression or a string specifying a regular expression.

Some examples:

```
;;; Replace occurrences of "Cotton" with "Jin".
(regexp-substitute/global #f (rx "Cotton") s
  'pre "Jin" 'post)

;;; mm/dd/yy -> dd/mm/yy date conversion.
(regexp-substitute/global #f (rx (submatch (+ digit)) "/" ; 1 = M
  (submatch (+ digit)) "/" ; 2 = D
  (submatch (+ digit))) ; 3 = Y
  s ; Source string
  'pre 2 "/" 1 "/" 3 'post)

;;; "9/29/61" -> "Sep 29, 1961" date conversion.
(regexp-substitute/global #f (rx (submatch (+ digit)) "/" ; 1 = M
  (submatch (+ digit)) "/" ; 2 = D
  (submatch (+ digit))) ; 3 = Y
  s ; Source string
  'pre
  ;; Sleazy converter -- ignores "year 2000" issue,
  ;; and blows up if month is out of range.
  (lambda (m)
    (let ((mon (vector-ref #'("Jan" "Feb" "Mar" "Apr" "May" "Jun"
      "Jul" "Aug" "Sep" "Oct" "Nov" "Dec")
      (- (string->number (match:substring m 1)) 1)))
      (day (match:substring m 2))
      (year (match:substring m 3)))
      (string-append mon " " day " ", 19" year)))
    'post)

;;; Remove potentially offensive substrings from string S.
(define (kill-matches re s)
  (regexp-substitute/global #f re s 'pre 'post))

(kill-matches (rx (| "Windows" "tcl" "Intel")) s) ; Protect the children.
```

(regexp-fold *re kons knil s [finish start]*)  $\longrightarrow$  *object* procedure

The following definition is a bit unwieldy, but the intuition is simple: this procedure uses the regexp *re* to divide up string *s* into non-matching/matching chunks, and then “folds” the procedure *kons* across this sequence of chunks. It is useful when you wish to operate on a string in sub-units defined by some regular expression, as are the related `regexp-fold-right` and `regexp-for-each` procedures.

Search from *start* (defaulting to 0) for a match to *re*; call this match *m*. Let *i* be the index of the end of the match (that is, `(match:end m 0)`). Loop as follows:

`(regexp-fold re kons (kons start m knil) s finish i)`

If there is no match, return instead

`(finish start knil)`

*Finish* defaults to `(lambda (i knil) knil)`.

In other words, we divide up *s* into a sequence of non-matching/matching chunks:

$$NM_1 M_1 NM_1 M_2 \dots NM_{k-1} M_{k-1} NM_k$$

where  $NM_1$  is the initial part of *s* that isn’t matched by the regexp *re*,  $M_1$  is the first match,  $NM_2$  is the following part of *s* that isn’t matched,  $M_2$  is the second match, and so forth— $NM_k$  is the final non-matching chunk of *s*. We apply *kons* from left to right to build up a result, passing it one non-matching/matching chunk each time: on an application `(kons i m knil)`, the non-matching chunk goes from *i* to `(match:begin m 0)`, and the following matching chunk goes from `(match:begin m 0)` to `(match:end m 0)`. The last non-matching chunk  $NM_k$  is processed by *k*. So the computation we perform is

$$(final\ Q\ (kons\ j_k\ M_k\ \dots\ (kons\ J_1\ M_1\ knil)\ \dots))$$

where  $J_i$  is the index of the start of  $NM_i$ ,  $M_i$  is a match value describing  $M_i$ , and *Q* is the index of the beginning of  $NM_k$ .

Hint: The `let-match` macro is frequently useful for operating on the match value *M* passed to the *kons* function.

`(regexp-fold-right re kons knil s [finish start])`  $\longrightarrow$  *object*      procedure

The right-to-left variant of `regexp-fold`.

This procedure repeatedly matches regexp *re* across string *s*. This divides *s* up into a sequence of matching/non-matching chunks:

$$NM_1 M_1 NM_1 M_2 \dots NM_{k-1} M_{k-1} NM_k$$

where  $NM_1$  is the initial part of *s* that isn’t matched by the regexp *re*,  $M_1$  is the first match,  $NM_2$  is the following part of *s* that isn’t matched,  $M_2$

is the second match, and so forth— $NM_k$  is the final non-matching chunk of  $s$ . We apply *kons* from right to left to build up a result, passing it one non-matching/matching chunk each time:

(*final*  $Q$  (*kons*  $M_1 j_1 \dots (kons M_k J_k knil) \dots$ ))

where  $MTCH_i$  is a match value describing  $M_i$ ,  $J_i$  is the index of the end of  $NM_i$  (or, equivalently, the beginning of  $M_{i+1}$ ), and  $Q$  is the index of the beginning of  $M_1$ . In other words, *KONS* is passed a match, an index describing the following non-matching text, and the value produced by folding the following text. The *FINAL* function “polishes off” the fold operation by handling the initial chunk of non-matching text ( $NM_0$ , above). *FINISH* defaults to  $(\lambda (i knil) knil)$

Example: To pick out all the matches to *re* in *s*, say

```
(regexp-fold-right re
  (lambda (m i lis)
    (cons (match:substring m 0) lis))
  '() s)
```

Hint: The *let-match* macro is frequently useful for operating on the match value *m* passed to the *kons* function.

(*regexp-for-each re proc s [start]*)  $\longrightarrow$  *undefined* procedure

Repeatedly match regexp *re* against string *s*. Apply *proc* to each match that is produced. Matches do not overlap.

Hint: The *let-match* macro is frequently useful for operating on the match value *m* passed to *varproc*.

(*let-match match-exp mvars body ...*)  $\longrightarrow$  *object* Syntax

(*if-match match-exp mvars on-match no-match*)  $\longrightarrow$  *object* Syntax

*Mvars* is a list of vars that is bound to the match and submatches of the string; *#F* is allowed as a don’t-care element. For example,

```
(let-match (regexp-search date s) (whole-date month day year)
  ... body ...)
```

matches the regexp against string *s*, then evaluates the body of the *let-match* in a scope where *whole-date* is bound to the matched string, and *month*, *day* and *year* are bound to the first, second and third submatches.

*if-match* is similar, but if the match expression is false, then the *no-match* expression is evaluated; this would be an error in *let-match*.

(*match-cond clause ...*)  $\longrightarrow$  *object* Syntax

This macro allows one to conditionally attempt a sequence of pattern matches, interspersed with other, general conditional tests. There are four kinds of `match-cond` clause, one introducing a pattern match, and the other three simply being regular `cond`-style clauses, marked by the `test` and `else` keywords:

```
(match-cond (match-exp match-vars body ...) ; As in if-match
            (test exp body ...)           ; As in cond
            (test exp => proc)              ; As in cond
            (else body ...))              ; As in cond
```

```
(flush-submatches re)  → re           procedure
(uncase re)            → re           procedure
(simplify-regexp re)   → re           procedure
(uncase-char-set cset) → re           procedure
(uncase-string str)    → re           procedure
```

These functions map regexps and char sets to other regexps. `flush-submatches` returns a regexp which matches exactly what its argument matches, but contains no submatches.

`uncase` returns a regexp that matches any case-permutation of its argument regexp.

`simplify-regexp` applies the simplifier to its argument. This is done automatically when compiling regular expressions, so this is only useful for programmers that are directly examining the ADT value with lower-level accessors.

`uncase-char-set` maps a char set to a regular expression that matches any character from that set, regardless of case. Similarly, `uncase-string` returns a regexp that matches any case-permutation of the string. For example, `(uncase-string "Knight")` returns the same value that `(rx ("kK") ("nN") ("iI") ("gG") ("hH") ("tT"))` or `(rx (w/nocase "Knight"))`.

```
(sre->regexp sre)  → re           procedure
(regexp->sre re)   → sre          procedure
```

These are the SRE parser and unparser. That is, `sre->regexp` maps an SRE to a regexp value, and `regexp->sre` does the inverse. The latter function can be useful for printing out regexps in a readable format.

```

(sre->regexp '(: "Olin " (? "G. ") "Shivers")) => regexp
(define re (re-seq (re-string "Pete ")
                  (re-repeat 1 #f (re-string "Sz"))
                  (re-string "ilagyi")))
(regexp->sre (re-repeat 0 1 re))
=> '(? "Pete" (+ "Sz") "ilagyi")

```

```

(posix-string->regexp string)  → re                                procedure
(regexp->posix-string re)      → [string syntax-level paren-count submatches-vector]  procedure

```

These two functions are the Posix notation parser and unparser. That is, `posix-string->regexp` maps a Posix-notation regular expression, such as `"g(ee|oo)se"`, to a regexp value, and `regexp->posix-string` does the inverse.

You can use these tools to map between scsh regexps and Posix regexp strings, which can be useful if you want to do conversion between SRE's and Posix form. For example, you can write a particularly complex regexp in SRE form, or compute it using the ADT constructors, then convert to Posix form, print it out, cut and paste it into a C or emacs lisp program. Or you can import an old regexp from some other program, parse it into an ADT value, render it to an SRE, print it out, then cut and paste it into a scsh program.

Note:

- The string parser doesn't handle the exotica of character class names such as `[[:alnum:]]`; the current implementation was written in in three hours.

## 6.6 The regexp ADT

The following functions may be used to construct and examine scsh's regexp abstract data type. They are in the following Scheme 48 packages: `re-adt-lib` `re-lib` `scsh`

Each basic class of regexp has a predicate, a basic constructor, a "smart" constructor that performs limited "peephole" optimisation on its arguments, and a set of accessors. The `...:tsm` accessor returns the total number of submatches contained in the regular expression.

<code>(re-seq? x)</code>	→ <i>boolean</i>	Type predicate
<code>(make-re-seq re-list)</code>	→ <i>re</i>	Basic constructor
<code>(re-seq re-list)</code>	→ <i>re</i>	Smart constructor
<code>(re-seq:elts re)</code>	→ <i>re-list</i>	Accessor

<code>(re-seq:tsm re)</code>	$\longrightarrow$ <i>integer</i>	Accessor
<code>(re-choice? x)</code>	$\longrightarrow$ <i>boolean</i>	Type predicate
<code>(make-re-choice re-list)</code>	$\longrightarrow$ <i>re</i>	Basic constructor
<code>(re-choice re-list)</code>	$\longrightarrow$ <i>re</i>	Smart constructor
<code>(re-choice:elts re)</code>	$\longrightarrow$ <i>re-list</i>	Accessor
<code>(re-choice:tsm re)</code>	$\longrightarrow$ <i>integer</i>	Accessor
<code>(re-repeat? x)</code>	$\longrightarrow$ <i>boolean</i>	Type predicate
<code>(make-re-repeat from to body)</code>	$\longrightarrow$ <i>re</i>	Accessor
<code>(re-repeat:from re)</code>	$\longrightarrow$ <i>integer</i>	Accessor
<code>(re-repeat:to re)</code>	$\longrightarrow$ <i>integer</i>	Accessor
<code>(re-repeat:tsm re)</code>	$\longrightarrow$ <i>integer</i>	Accessor
<code>(re-submatch? x)</code>	$\longrightarrow$ <i>boolean</i>	Type predicate
<code>(make-re-submatch body [pre-dsm post-dsm])</code>	$\longrightarrow$ <i>re</i>	Accessor
<code>(re-submatch:pre-dsm re)</code>	$\longrightarrow$ <i>integer</i>	Accessor
<code>(re-submatch:post-dsm re)</code>	$\longrightarrow$ <i>integer</i>	Accessor
<code>(re-submatch:tsm re)</code>	$\longrightarrow$ <i>integer</i>	Accessor
<code>(re-string? x)</code>	$\longrightarrow$ <i>boolean</i>	Type predicate
<code>(make-re-string chars)</code>	$\longrightarrow$ <i>re</i>	Basic constructor
<code>(re-string chars)</code>	$\longrightarrow$ <i>re</i>	Basic constructor
<code>(re-string:chars re)</code>	$\longrightarrow$ <i>string</i>	Accessor
<code>(re-char-set? x)</code>	$\longrightarrow$ <i>boolean</i>	Type predicate
<code>(make-re-char-set cset)</code>	$\longrightarrow$ <i>re</i>	Basic constructor
<code>(re-char-set cset)</code>	$\longrightarrow$ <i>re</i>	Basic constructor
<code>(re-char-set:cset re)</code>	$\longrightarrow$ <i>char-set</i>	Accessor
<code>(re-dsm? x)</code>	$\longrightarrow$ <i>boolean</i>	Type predicate
<code>(make-re-dsm body pre-dsm post-dsm)</code>	$\longrightarrow$ <i>re</i>	Basic constructor
<code>(re-dsm body pre-dsm post-dsm)</code>	$\longrightarrow$ <i>re</i>	Smart constructor
<code>(re-dsm:body re)</code>	$\longrightarrow$ <i>re</i>	Accessor
<code>(re-dsm:pre-dsm re)</code>	$\longrightarrow$ <i>integer</i>	Accessor
<code>(re-dsm:post-dsm re)</code>	$\longrightarrow$ <i>integer</i>	Accessor
<code>(re-dsm:tsm re)</code>	$\longrightarrow$ <i>integer</i>	Accessor
<code>re-bos</code>		regexp
<code>re-eos</code>		regexp
<code>re-bol</code>		regexp
<code>re-eol</code>		regexp

These variables are bound to the primitive anchor regexps.

<code>(re-bos? object)</code>	$\longrightarrow$ <i>boolean</i>	procedure
<code>(re-eos? object)</code>	$\longrightarrow$ <i>boolean</i>	procedure
<code>(re-bol? object)</code>	$\longrightarrow$ <i>boolean</i>	procedure
<code>(re-eol? object)</code>	$\longrightarrow$ <i>boolean</i>	procedure



These predicates recognise the associated primitive anchor regexp.

<code>re-trivial</code>	regexp
<code>(re-trivial? re) → boolean</code>	procedure

The variable `re-trivial` is bound to a regular expression that matches the empty string (corresponding to the SRE `""` or `(:)`); it is recognised by the associated predicate. Note that the predicate is only guaranteed to recognise this particular trivial regexp; other trivial regexps built using other constructors may or may not produce a true value.

<code>re-empty</code>	regexp
<code>(re-empty? re) → boolean</code>	procedure

The variable `re-empty` is bound to a regular expression that never matches (corresponding to the SRE `()`); it is recognised by the associated predicate. Note that the predicate is only guaranteed to recognise this particular empty regexp; other empty regexps built using other constructors may or may not produce a true value.

<code>re-any</code>	regexp
<code>(re-any? re) → boolean</code>	procedure

The variable `re-any` is bound to a regular expression that matches any character (corresponding to the SRE `any`); it is recognised by the associated predicate. Note that the predicate is only guaranteed to recognise this particular any-character regexp value; other any-character regexps built using other constructors may or may not produce a true value.

<code>re-nonl</code>	regexp
----------------------	--------

The variable `re-nonl` is bound to a regular expression that matches any non-newline character (corresponding to the SRE `(~ #\newline)`).

<code>(regexp? object) → boolean</code>	procedure
---	-----------

Is the object a regexp?

<code>(re-tsm re) → integer</code>	procedure
------------------------------------	-----------

Return the total number of submatches contained in the regexp.

<code>(clean-up-cres) → undefined</code>	procedure
--	-----------

The current `scsh` implementation should call this function periodically to release C-heap storage associated with compiled regexps. Hopefully, this procedure will be removed at a later date.

## 6.7 Syntax-hacking tools

The Scheme 48 package `sre-syntax-tools` exports several tools for macro writers that want to use SREs in their macros. In the functions defined below, *compare* and *rename* parameters are as passed to Clinger-Rees explicit-renaming low-level macros.

`(if-sre-form form conseq-form alt-form) → form` Syntax

If *form* is a legal SRE, this is equivalent to the expression *conseq-form*, otherwise it expands to *alt-form*.

This is useful for high-level macro authors who want to write a macro where one field in the macro can be an SRE or possibly something else. E.g., we might have a conditional form wherein if the test part of one arm is an SRE, it expands to a regexp match on some implied value, otherwise the form is evaluated as a boolean Scheme expression. For example, a conditional macro might expand into code containing the following form, which in turn would have one of two possible expansions:

```
(if-sre-form test-exp          ; If TEST-EXP is SRE,
  (regexp-search? (rx test-exp) line) ; match it w/the line,
  test-exp)                       ; otw it's a text exp.
```

`(sre-form? form rename compare) → boolean` procedure

This procedure is for low-level macros doing things equivalent to `if-sre-form`. It returns true if the form is a legal SRE.

Note that neither `sre-form` nor `if-sre-form` does a deep recursion over the form in the case where the form is a list. They simply check the car of the form for one of the legal SRE keywords.

`(parse-sre sre-form compare rename) → re` procedure

`(parse-sres sre-forms compare rename) → re` procedure

Parse *sre-form* into an ADT. Note that if the SRE is dynamic—contains *,exp* or *,@exp* forms, or has repeat operators whose from/to counts are not constants—then the returned ADT will have *Scheme expressions* in the corresponding slots of the regexp records instead of the corresponding integer, char-set, or regexp. In other words, we use the ADT as its own AST. It's called a "hack."

`parse-sres` parses a list of SRE forms that comprise an implicit sequence.

`(regexp->scheme re rename) → Scheme-expression` procedure

Returns a Scheme expression that will construct the regexp *re* using ADT constructors such as `make-re-sequence`, `make-re-repeat`, and so forth.

If the regexp is static, it will be simplified and pre-translated to a Posix string as well, which will be part of the constructed regexp value.

`(static-regexp? re) → boolean`

Is the regexp a static one?

procedure

## Chapter 7

# Reading delimited strings

Scsh provides a set of procedures that read delimited strings from input ports. There are procedures to read a single line of text (terminated by a newline character), a single paragraph (terminated by a blank line), and general delimited strings (terminated by a character belonging to an arbitrary character set).

These procedures can be applied to any Scheme input port. However, the scsh virtual machine has native-code support for performing delimited reads on Unix ports, and these input operations should be particularly fast—much faster than doing the equivalent character-at-a-time operation from Scheme code.

All of the delimited input operations described below take a `handle-delim` parameter, which determines what the procedure does with the terminating delimiter character. There are four possible choices for a `handle-delim` parameter:

handle-delim	Meaning
'trim	Ignore delimiter character.
'peek	Leave delimiter character in input stream.
'concat	Append delimiter character to returned value.
'split	Return delimiter as second value.

The first case, `'trim`, is the standard default for all the routines described in this section. The last three cases allow the programmer to distinguish between strings that are terminated by a delimiter character, and strings that are terminated by an end-of-file.

`(read-line [port handle-newline])`  $\longrightarrow$  *string or eof-object*      procedure  
Reads and returns one line of text; on eof, returns the eof object. A line is terminated by newline or eof.

*handle-newline* determines what `read-line` does with the newline or EOF that terminates the line; it takes the general set of values described for the general *handle-delim* case above, and defaults to `'trim` (discard the newline). Using this argument allows one to tell whether or not the last line of input in a file is newline terminated.

`(read-paragraph [port handle-delim])`  $\longrightarrow$  *string or eof*      procedure

This procedure skips blank lines, then reads text from a port until a blank line or eof is found. A “blank line” is a (possibly empty) line composed only of white space. The *handle-delim* parameter determines how the terminating blank line is handled. It is described above, and defaults to `'trim`. The `'peek` option is not available.

The following procedures read in strings from ports delimited by characters belonging to a specific set. See section 5.5 for information on character set manipulation.

`(read-delimited char-set [port handle-delim])`  $\longrightarrow$  *string or eof*      procedure

Read until we encounter one of the chars in *char-set* or eof. The *handle-delim* parameter determines how the terminating character is handled. It is described above, and defaults to `'trim`.

The *char-set* argument may be a charset, a string, a character, or a character predicate; it is coerced to a charset.

`(read-delimited! char-set buf [port handle-delim start end])`  $\longrightarrow$  *nchars or eof or #f*      procedure

A side-effecting variant of `read-delimited`.

The data is written into the string *buf* at the indices in the half-open interval  $[start, end)$ ; the default interval is the whole string:  $start = 0$  and  $end = (\text{string-length } buf)$ . The values of *start* and *end* must specify a well-defined interval in *str*, i.e.,  $0 \leq start \leq end \leq (\text{string-length } buf)$ .

It returns *nbytes*, the number of bytes read. If the buffer filled up without a delimiter character being found, `#f` is returned. If the port is at eof when the read starts, the eof object is returned.

If an integer is returned (i.e., the read is successfully terminated by reading a delimiter character), then the *handle-delim* parameter determines how the terminating character is handled. It is described above, and defaults to `'trim`.

`(%read-delimited! char-set buf gobble? [port start end])`  $\longrightarrow$  *[char-or-eof-or-#f integer]*      procedure

This low-level delimited reader uses an alternate interface. It returns two values: *terminator* and *num-read*.

**terminator** A value describing why the read was terminated:

Character or eof-object	$\Rightarrow$	Read terminated by this value.
#f	$\Rightarrow$	Filled buffer without finding a delimiter.

**num-read** Number of characters read into *buf*.

If the read is successfully terminated by reading a delimiter character, then the *gobble?* parameter determines what to do with the terminating character. If true, the character is removed from the input stream; if false, the character is left in the input stream where a subsequent read operation will retrieve it. In either case, the character is also the first value returned by the procedure call.

(skip-char-set *skip-chars* [*port*])  $\longrightarrow$  *integer* procedure

Skip characters occurring in the set *skip-chars*; return the number of characters skipped. The *skip-chars* argument may be a charset, a string, a character, or a character predicate; it is coerced to a charset.

## Chapter 8

# Awk, record I/O, and field parsing

Unix programs frequently process streams of records, where each record is delimited by a newline, and records are broken into fields with other delimiters (for example, the colon character in `/etc/passwd`). Scsh has procedures that allow the programmer to easily do this kind of processing. Scsh's field parsers can also be used to parse other kinds of delimited strings, such as colon-separated `$PATH` lists. These routines can be used with scsh's `awk` loop construct to conveniently perform pattern-directed computation over streams of records.

### 8.1 Record I/O and field parsing

The procedures in this section are used to read records from I/O streams and parse them into fields. A record is defined as text terminated by some delimiter (usually a newline). A record can be split into fields by using regular expressions in one of several ways: to *match* fields, to *separate* fields, or to *terminate* fields. The field parsers can be applied to arbitrary strings (one common use is splitting environment variables such as `$PATH` at colons into its component elements).

The general delimited-input procedures described in chapter 7 are also useful for reading simple records, such as single lines, paragraphs of text, or strings terminated by specific characters.

### 8.1.1 Reading records

(record-reader [*delims elide-delims? handle-delim*])  $\longrightarrow$  procedure procedure

Returns a procedure that reads records from a port. The procedure is invoked as follows:

(reader [*port*])  $\longrightarrow$  string or eof

A record is a sequence of characters terminated by one of the characters in *delims* or eof. If *elide-delims?* is true, then a contiguous sequence of delimiter chars are taken as a single record delimiter. If *elide-delims?* is false, then a delimiter char coming immediately after a delimiter char produces an empty-string record. The reader consumes the delimiting char(s) before returning from a read.

The *delims* set defaults to the set {newline}. It may be a charset, string, character, or character predicate, and is coerced to a charset. The *elide-delims?* flag defaults to #f.

The *handle-delim* argument controls what is done with the record's terminating delimiter.

- 'trim Delimiters are trimmed. (The default)
- 'split Reader returns delimiter string as a second argument. If record is terminated by EOF, then the eof object is returned as this second argument.
- 'concat The record and its delimiter are returned as a single string.

The reader procedure returned takes one optional argument, the port from which to read, which defaults to the current input port. It returns a string or eof.

### 8.1.2 Parsing fields

(field-splitter [*field num-fields*])  $\longrightarrow$  procedure procedure  
(infix-splitter [*delim num-fields handle-delim*])  $\longrightarrow$  procedure procedure  
(suffix-splitter [*delim num-fields handle-delim*])  $\longrightarrow$  procedure procedure  
(sloppy-suffix-splitter [*delim num-fields handle-delim*])  $\longrightarrow$  procedure procedure

These functions return a parser function that can be used as follows:

(parser string [*start*])  $\longrightarrow$  string-list

The returned parsers split strings into fields defined by regular expressions. You can parse by specifying a pattern that *separates* fields, a pattern that *terminates* fields, or a pattern that *matches* fields:



Procedure	Pattern
<code>field-splitter</code>	matches fields
<code>infix-splitter</code>	separates fields
<code>suffix-splitter</code>	terminates fields
<code>sloppy-suffix-splitter</code>	terminates fields

These parser generators are controlled by a range of options, so that you can precisely specify what kind of parsing you want. However, these options default to reasonable values for general use.

Defaults:

```

delim          = (rx (| (+ white) eos)) (suffix delimiter: white space or eos)
                  (rx (+ white))        (infix delimiter: white space)
field          = (rx (+ (~ white)))    (non-white-space)
num-fields     = #f                    (as many fields as possible)
handle-delim  = 'trim                  (discard delimiter chars)

```

...which means: break the string at white space, discarding the white space, and parse as many fields as possible.

The *delim* parameter is a regular expression matching the text that occurs between fields. See chapter 6 for information on regular expressions, and the `rx` form used to specify them. In the separator case, it defaults to a pattern matching white space; in the terminator case, it defaults to white space or end-of-string.

The *field* parameter is a regular expression used to match fields. It defaults to non-white-space.

The *delim* patterns may also be given as a string, character, or char-set, which are coerced to regular expressions. So the following expressions are all equivalent, each producing a function that splits strings apart at colons:

```

(infix-splitter (rx ":"))
(infix-splitter ":")
(infix-splitter #\:)
(infix-splitter (char-set #\:))

```

The boolean *handle-delim* determines what to do with delimiters.

```

'trim    Delimiters are thrown away after parsing. (default)
'concat  Delimiters are appended to the field preceding them.
'split   Delimiters are returned as separate elements in the field list.

```

The *num-fields* argument used to create the parser specifies how many fields to parse. If `#f` (the default), the procedure parses them all. If a positive integer *n*, exactly that many fields are parsed; it is an error if there are more or fewer than *n* fields in the record. If *num-fields* is a negative integer or zero, then  $|n|$  fields are parsed, and the remainder of the string

is returned in the last element of the field list; it is an error if fewer than  $|n|$  fields can be parsed.

The field parser produced is a procedure that can be employed as follows:

$(\text{parse string } [start]) \implies \text{string-list}$

The optional *start* argument (default 0) specifies where in the string to begin the parse. It is an error if  $start > (\text{string-length string})$ .

The parsers returned by the four parser generators implement different kinds of field parsing:

**field-splitter** The regular expression specifies the actual field.

**suffix-splitter** Delimiters are interpreted as element *terminators*. If vertical-bar is the the delimiter, then the string "" is the empty record (), "foo|" produces a one-field record ("foo"), and "foo" is an error.

The syntax of suffix-delimited records is:

$\langle \text{record} \rangle ::= \text{""} \quad (\text{Empty record})$   
 $\quad \quad \quad | \quad \langle \text{element} \rangle \langle \text{delim} \rangle \langle \text{record} \rangle$

It is an error if a non-empty record does not end with a delimiter. To make the last delimiter optional, make sure the delimiter regexp matches the end-of-string (sre eos).

**infix-splitter** Delimiters are interpreted as element *separators*. If comma is the delimiter, then the string "foo," produces a two-field record ("foo" "").

The syntax of infix-delimited records is:

$\langle \text{record} \rangle ::= \text{""} \quad (\text{Forced to be empty record})$   
 $\quad \quad \quad | \quad \langle \text{real-infix-record} \rangle$

$\langle \text{real-infix-record} \rangle ::= \langle \text{element} \rangle \langle \text{delim} \rangle \langle \text{real-infix-record} \rangle$   
 $\quad \quad \quad | \quad \langle \text{element} \rangle$

Note that separator semantics doesn't really allow for empty records—the straightforward grammar (i.e.,  $\langle \text{real-infix-record} \rangle$ ) parses an empty string as a singleton list whose one field is the empty string, (""), not as the empty record (). This is unfortunate, since it means that infix string parsing doesn't make string-append and append isomorphic. For example,

((infix-splitter ":") (string-append x ":" y))

doesn't always equal

(append ((infix-splitter ":") x)  
((infix-splitter ":") y))

Record	: suffix	: \$ suffix	: infix	non-: field
" "	()	()	()	()
": "	(" ")	(" ")	(" " " ")	()
"foo: "	("foo")	("foo")	("foo" " ")	("foo")
":foo"	<i>error</i>	(" " "foo")	(" " "foo")	("foo")
"foo:bar"	<i>error</i>	("foo" "bar")	("foo" "bar")	("foo" "bar")

Figure 8.1: Using different grammars to split records into fields.

It fails when  $x$  or  $y$  are the empty string. Terminator semantics *does* preserve a similar isomorphism.

However, separator semantics is frequently what other Unix software uses, so to parse their strings, we need to use it. For example, Unix \$PATH lists have separator semantics. The path list `"/bin: "` is broken up into `("/bin" " ")`, not `("/bin")`. Comma-separated lists should also be parsed this way.

**sloppy-suffix** The same as the **suffix** case, except that the parser will skip an initial delimiter string if the string begins with one instead of parsing an initial empty field. This can be used, for example, to field-split a sequence of English text at white-space boundaries, where the string may begin or end with white space, by using regex

```
(rx (| (+ white) eos))
```

(But you would be better off using `field-splitter` in this case.)

Figure 8.1 shows how the different parser grammars split apart the same strings. Having to choose between the different grammars requires you to decide what you want, but at least you can be precise about what you are parsing. Take fifteen seconds and think it out. Say what you mean; mean what you say.

```
(join-strings string-list [delimiter grammar])  →  string      procedure
```

This procedure is a simple unparser—it pastes strings together using the delimiter string.

The *grammar* argument is one of the symbols `infix` (the default) or `suffix`; it determines whether the delimiter string is used as a separator or as a terminator.

The delimiter is the string used to delimit elements; it defaults to a single space `" "`.

Example:

```
(join-strings '("foo" "bar" "baz") ":")
⇒ "foo:bar:baz"
```

### 8.1.3 Field readers

(field-reader [*field-parser* *rec-reader*]) → *procedure*                      procedure

This utility returns a procedure that reads records with field structure from a port. The reader's interface is designed to make it useful in the `awk` loop macro (section 8.2). The reader is used as follows:

```
(reader [port]) ⇒ [raw-record parsed-record] or [eof()]
```

When the reader is applied to an input port (default: the current input port), it reads a record using *rec-reader*. If this record isn't the eof object, it is parsed with *field-parser*. These two values—the record, and its parsed representation—are returned as multiple values from the reader.

When called at eof, the reader returns [eof-object ()].

Although the record reader typically returns a string, and the field-parser typically takes a string argument, this is not required. The record reader can produce, and the field-parser consume, values of any type. However, the empty list returned as the parsed value on eof is hardwired into the field reader.

For example, if port `p` is open on `/etc/passwd`, then

```
((field-reader (infix-splitter ":" 7)) p)
```

returns two values:

```
"dalbertz:mx3Uaqq0:107:22:David Albertz:/users/dalbertz:/bin/csh"
("dalbertz" "mx3Uaqq0" "107" "22" "David Albertz" "/users/dalbertz"
 "/bin/csh")
```

The *field-parser* defaults to the value of (field-splitter), a parser that picks out sequences of non-white-space strings.

The *rec-reader* defaults to read-line.

Figure 8.2 shows field-reader being used to read different kinds of Unix records.

### 8.1.4 Forward-progress guarantees and empty-string matches

A loop that pulls text off a string by repeatedly matching a regexp against that string can conceivably get stuck in an infinite loop if the regexp matches the empty string. For example, the SREs `bos`, `eos`, `(* any)`, and `(| "foo" (* ( "f")))` can all match the empty string.

```

;;; /etc/passwd reader
(field-reader (infix-splitter ":" 7))
  ; wandy:3xuncWdpKhR.:73:22:Wandy Saetan:/usr/wandy:/bin/csh

;;; Two ls -l output readers
(field-reader (infix-splitter (rx (+ white)) 8))
(field-reader (infix-splitter (rx (+ white)) -7))
  ; -rw-r--r--  1 shivers    22880 Sep 24 12:45 scsh.scm

;;; Internet hostname reader
(field-reader (field-splitter (rx (+ (~ "."))))))
  ; stat.sinica.edu.tw

;;; Internet IP address reader
(field-reader (field-splitter (rx (+ (~ "."))) 4))
  ; 18.24.0.241

;;; Line of integers
(let ((parser (field-splitter (rx (? ("+-")) (+ digit)))))
  (field-reader (lambda (s) (map string->number (parser s))))
  ; 18 24 0 241

;;; Same as above.
(let ((reader (field-reader (field-splitter (rx (? ("+-"))
                                              (+ digit))))))
  (lambda (maybe-port) (map string->number (apply reader maybe-port))))
  ; Yale beat harvard 26 to 7.

```

Figure 8.2: Some examples of field-reader

The routines in this package that iterate through strings with regular expressions are careful to handle this empty-string case. If a regexp matches the empty string, the next search starts, not from the end of the match (which in the empty string case is also the beginning—that’s the problem), but from the next character over. This is the correct behaviour. Regexp match the longest possible string at a given location, so if the regexp matched the empty string at location  $i$ , then it is guaranteed it could not have matched a longer pattern starting with character  $i$ . So we can safely begin our search for the next match at char  $i + 1$ .

With this provision, every iteration through the loop makes some forward progress, and the loop is guaranteed to terminate.

This has the effect you want with field parsing. For example, if you split a string with the empty pattern, you will explode the string into its individual characters:

```
((suffix-splitter (rx)) "foo") => (" " "f" "o" "o")
```

However, even though this boundary case is handled correctly, we don’t recommend using it. Say what you mean—just use a field splitter:

```
((field-splitter (rx any)) "foo") => ("f" "o" "o")
```

Or, more efficiently,

```
((λ (s) (map string (string->list s))) "foo")
```

### 8.1.5 Reader limitations

Since all of the readers in this package require the ability to peek ahead one char in the input stream, they cannot be applied to raw integer file descriptors, only Scheme input ports. This is because Unix doesn’t support peeking ahead into input streams.

## 8.2 Awk

Scsh provides a loop macro and a set of field parsers that can be used to perform text processing very similar to the Awk programming language. The basic functionality of Awk is factored in scsh into its component parts. The control structure is provided by the awk loop macro; the text I/O and parsers are provided by the field-reader subroutine library (section 8.1). This factoring allows the programmer to compose the basic loop structure with any parser or input mechanism at all. If the parsers provided by the field-reader package are insufficient, the programmer can write a custom parser in Scheme and use it with equal ease in the awk framework.

Awk-in-scheme is given by a loop macro called `awk`. It looks like this:

```
(awk <next-record> <record&field-vars>
    [<counter>] <state-var-decls>
    <clause1> ...)
```

The body of the loop is a series of clauses, each one representing a kind of condition/action pair. The loop repeatedly reads a record, and then executes each clause whose condition is satisfied by the record.

Here's an example that reads lines from port `p` and prints the line number and line of every line containing the string "Church-Rosser":

```
(awk (read-line) (ln) lineno ()
    ("Church-Rosser" (format #t "~d: ~s~%" lineno ln)))
```

This example has just one clause in the loop body, the one that tests for matches against the regular expression "Church-Rosser".

The `<next-record>` form is an expression that is evaluated each time through the loop to produce a record to process. This expression can return multiple values; these values are bound to the variables given in the `<record&field-vars>` list of variables. The first value returned is assumed to be the record; when it is the end-of-file object, the loop terminates.

For example, let's suppose we want to read items from `/etc/password`, and we use the `field-reader` procedure to define a record parser for `/etc/passwd` entries:

```
(define read-passwd (field-reader (infix-splitter ":" 7)))
```

binds `read-passwd` to a procedure that reads in a line of text when it is called, and splits the text at colons. It returns two values: the entire line read, and a seven-element list of the split-out fields. (See section 8.1 for more on `field-reader` and `infix-splitter`.)

So if the `<next-record>` form in an `awk` expression is `(read-passwd)`, then `<record&field-vars>` must be a list of two variables, *e.g.*,

```
(record field-vec)
```

since `read-passwd` returns two values.

Note that `awk` allows us to use *any* record reader we want in the loop, returning whatever number of values we like. These values don't have to be strings or string lists. The only requirement is that the record reader return the eof object as its first value when the loop should terminate.

The `awk` loop allows the programmer to have loop variables. These are declared and initialised by the `<state-var-decls>` form, a

```
((var init-exp) (var init-exp) ...)
```

list rather like the `let` form. Whenever a clause in the loop body executes, it evaluates to as many values as there are state variables, updating them.

The optional  $\langle counter \rangle$  variable is an iteration counter. It is bound to 0 when the loop starts. The counter is incremented each time a non-eof record is read.

There are several kinds of loop clause. When evaluating the body of the loop, `awk` evaluates *all* the clauses sequentially. Unlike `cond`, it does not stop after the first clause is satisfied; it checks them all.

- $(test\ body_1\ body_2\ \dots)$

If *test* is true, execute the body forms. The last body form is the value of the clause. The test and body forms are evaluated in the scope of the record and state variables.

The *test* form can be one of:

<i>integer</i> :	The test is true for that iteration of the loop. The first iteration is #1.
<i>sre</i> :	A regular expression, in SRE notation (see chapter 6) can be used as a test. The test is successful if the pattern matches the record. In particular, note that any string is an SRE.
$(when\ expr)$ :	The body of a <code>when</code> test is evaluated as a Scheme boolean expression in the inner scope of the <code>awk</code> form.
<i>expr</i> :	If the form is none of the above, it is treated as a Scheme expression—in practice, the <code>when</code> keyword is only needed in cases where SRE/Scheme expression ambiguity might occur.

- $(range\ start-test\ stop-test\ body_1\ \dots)$   
 $(:range\ start-test\ stop-test\ body_1\ \dots)$   
 $(range:\ start-test\ stop-test\ body_1\ \dots)$   
 $(:range:\ start-test\ stop-test\ body_1\ \dots)$

These clauses become activated when *start-test* is true; they stay active on all further iterations until *stop-test* is true.

So, to print out the first ten lines of a file, we use the clause:

```
(:range: 1 10 (display record))
```

The colons control whether or not the start and stop lines are processed by the clause. For example:



```

(range 1 5 ...)    Lines 2 3 4
(:range 1 5 ...)   Lines 1 2 3 4
(range: 1 5 ...)    Lines 2 3 4 5
(:range: 1 5 ...)   Lines 1 2 3 4 5

```

A line can trigger both tests, either simultaneously starting and stopping an active region, or simultaneously stopping one and starting a new one, so ranges can abut seamlessly.

- `(else body1 body2 ...)`  
If no other clause has executed since the top of the loop, or since the last else clause, this clause executes.
- `(test => exp)`  
If evaluating *test* produces a true value, apply *exp* to that value. If *test* is a regular expression, then *exp* is applied to the match data structure returned by the regexp match routine.
- `(after body1 ...)`  
This clause executes when the loop encounters EOF. The body forms execute in the scope of the state vars and the record-count var, if there are any. The value of the last body form is the value of the entire awk form.  
  
If there is no after clause, awk returns the loop's state variables as multiple values.

### 8.2.1 Examples

Here are some examples of awk being used to process various types of input stream.

```

(define $ list-ref)      ; Saves typing.

;;; Print out the name and home-directory of everyone in /etc/passwd:
(let ((read-passwd (field-reader (infix-splitter ":" 7))))
  (call-with-input-file "/etc/passwd"
    (lambda (port)
      (awk (read-passwd port) (record fields) ()
        (#t (format #t "~a's home directory is ~a%"
                    ($ fields 0)
                    ($ fields 5)))))))

```

```

;;; Print out the user-name and home-directory of everyone whose
;;; name begins with "S"
(let ((read-passwd (field-reader (infix-splitter ":" 7))))
  (call-with-input-file "/etc/passwd"
    (lambda (port)
      (awk (read-passwd port) (record fields) ()
        ( (: bos "S")
          (format #t "~a's home directory is ~a~%"
            ($ fields 0)
            ($ fields 5)))))))

;;; Read a series of integers from stdin. This expression evaluates
;;; to the number of positive numbers that were read. Note our
;;; "record-reader" is the standard Scheme READ procedure.
(awk (read) (i) ((npos 0))
  (> i 0) (+ npos 1)))

;;; Filter -- pass only lines containing my name.
(awk (read-line) (line) ()
  ("Olin" (display line) (newline)))

;;; Count the number of non-comment lines of code in my Scheme source.
(awk (read-line) (line) ((nlines 0))
  ( (: bos (* white) ";") nlines) ; A comment line.
  (else (+ nlines 1))) ; Not a comment line.

;;; Read numbers, counting the evens and odds.
(awk (read) (val) ((evens 0) (odds 0))
  (> val 0) (display "pos ") (values evens odds)) ; Tell me about
  (< val 0) (display "neg ") (values evens odds)) ; sign, too.
  (else (display "zero ") (values evens odds))

  ((even? val) (values (+ evens 1) odds))
  (else (values evens (+ odds 1))))

;;; Determine the max length of all the lines in the file.
(awk (read-line) (line) ((max-len 0))
  (#t (max max-len (string-length line))))

```

```

;;; (This could also be done with PORT-FOLD:)
(port-fold (current-input-port) read-line
  (lambda (line maxlen) (max (string-length line) maxlen))
  0)

;;; Print every line longer than 80 chars.
;;; Prefix each line with its line #.
(awk (read-line) (line) lineno ()
  ((> (string-length line) 80)
   (format #t "~d: ~s~%" lineno line)))

;;; Strip blank lines from input.
(awk (read-line) (line) ()
  ((~ white) (display line) (newline)))

;;; Sort the entries in /etc/passwd by login name.
(for-each (lambda (entry) (display (cdr entry)) (newline)) ; Out
  (sort (lambda (x y) (string<? (car x) (car y))) ; Sort
    (let ((read (field-reader (infix-splitter ":" 7)))) ; In
      (awk (read) (line fields) ((ans '()))
        (#t (cons (cons ($ fields 0) line) ans))))))

;;; Prefix line numbers to the input stream.
(awk (read-line) (line) lineno ()
  (#t (format #t "~d:\t~a~%" lineno line)))

```

### 8.3 Backwards compatibility

Previous scsh releases provided an awk form with a different syntax, designed around regular expressions written in Posix notation as strings, rather than SREs.

This form is still available in a separate module for old code. It'll be documented in the next release of this manual. Dig around in the sources for it.

## Chapter 9

# Concurrent system programming

The Scheme Shell provides the user with support for concurrent programming. The interface consists of several parts:

- The thread system
- Synchronization vehicles
- Process state abstractions

Whereas the user deals with threads and synchronization explicitly, the process state abstractions are built into the rest of the system, almost transparent for the user. Section 9.5 describes the interaction between process state and threads.

### 9.1 Threads

A thread can be thought of as a procedure that can run independently of and concurrent to the rest of the system. The calling procedure fires the thread up and forgets about it.

The current thread interface is completely taken from Scheme 48. This documentation is an extension of the file `doc/threads.txt`.

The thread structure is named `threads`, it has to be opened explicitly.

`(spawn thunk [name])`  $\longrightarrow$  `undefined` procedure

Create and schedule a new thread that will execute *thunk*, a procedure with no arguments. Note that Scsh's `spawn` does **not** return a reference to a thread object. The optional argument *name* is used when printing the thread.

The new thread will not inherit the values for the process state from its parent, see the procedure `fork-thread` in Section 9.5 for a way to create a thread with semantics similar to process forking.

`(relinquish-timeslice)`  $\longrightarrow$  *undefined* procedure

Let other threads run for a while.

`(sleep time)`  $\longrightarrow$  *undefined* procedure

Puts the current thread into sleep for *time* milliseconds. The time at which the thread is run again may be longer of course.

`(terminate-current-thread)`  $\longrightarrow$  *does-not-return* procedure

Kill the current thread.

Mainly for debugging purposes, there is also an interface to the internal representation of thread objects:

`(current-thread)`  $\longrightarrow$  *thread-object* procedure

Return the object to which the current thread internally corresponds. Note that this procedure is exported by the package `threads-internal` only.

`(thread? thing)`  $\longrightarrow$  *boolean* procedure

Returns true iff *thing* is a thread object.

`(thread-name thread)`  $\longrightarrow$  *name* procedure

*Name* corresponds to the second parameter that was given to `spawn` when *thread* was created.

`(thread-uid thread)`  $\longrightarrow$  *integer* procedure

Returns a unique identifier for the current thread.

## 9.2 Locks

Locks are a simple mean for mutual exclusion. They implement a concept commonly known as *semaphores*. Threads can obtain and release locks. If a thread tries to obtain a lock which is held by another thread, the first thread is blocked. To access the following procedures, you must open the structure `locks`.

`(make-lock)`  $\longrightarrow$  *lock* procedure

Creates a lock.

`(lock? thing)`  $\longrightarrow$  *boolean* procedure

Returns true iff *thing* is a lock.

`(obtain-lock lock)`  $\longrightarrow$  *undefined* procedure

Obtain *lock*. Causes the thread to block if the lock is held by a thread.

(maybe-obtain-lock *lock*)  $\longrightarrow$  *boolean* procedure  
 Tries to obtain *lock*, but returns false if the lock cannot be obtained.

(release-lock *lock*)  $\longrightarrow$  *boolean* procedure  
 Releases *lock*. Returns true if the lock immediately got a new owner, false otherwise.

(lock-owner-uid *lock*)  $\longrightarrow$  *integer* procedure  
 Returns the uid of the thread that currently holds *lock* or false if the lock is free.

### 9.3 Placeholders

Placeholders combine synchronization with value delivery. They can be thought of as special variables. After creation the value of the placeholder is undefined. If a thread tries to read the placeholders value this thread is blocked. Multiple threads are allowed to block on a single placeholder. They will continue running after another thread sets the value of the placeholder. Now all reading threads receive the value and continue executing. Setting a placeholder to two different values causes an error. The structure `placeholders` features the following procedures:

(make-placeholder)  $\longrightarrow$  *placeholder* procedure  
 Creates a new placeholder.

(placeholder? *thing*)  $\longrightarrow$  *boolean* procedure  
 Returns true iff *thing* is a placeholder.

(placeholder-set! *placeholder value*)  $\longrightarrow$  *undefined* procedure  
 Sets the placeholders value to *value*. If the placeholder is already set to a *different* value an exception is risen.

(placeholder-value *placeholder*)  $\longrightarrow$  *value* procedure  
 Returns the value of the placeholder. If the placeholder is yet unset, the current thread is blocked until another thread sets the value by means of `placeholder-set!`.

### 9.4 The event interface to interrupts

Scsh provides an synchronous interface to the asynchronous signals delivered by the operation system<sup>1</sup>. The key element in this system is an object called

<sup>1</sup>Olin's paper "Automatic management of operation-system resources" describes this system in detail.

*sigevent* which corresponds to the single occurrence of a signal. A *sigevent* has two fields: the Unix signal that occurred and a pointer to the *sigevent* that happened or will happen. That is, events are kept in a linked list in increasing-time order. Scsh's structure *sigevents* provides various procedures to access this list:

`(most-recent-sigevent)`  $\longrightarrow$  *sigevent* procedure

Returns the most recent *sigevent* — the head of the *sigevent* list.

`(sigevent? object)`  $\longrightarrow$  *boolean* procedure

The predicate for *sigevents*.

`(next-sigevent pre-event type)`  $\longrightarrow$  *event* procedure

Returns the next *sigevent* of type *type* after *sigevent pre-event*. If no such event exists, the procedure blocks.

`(next-sigevent-set pre-event set)`  $\longrightarrow$  *event* procedure

Returns the next *sigevent* whose type is in *set* after *pre-event*. If no such event exists, the procedure blocks.

`(next-sigevent/no-wait pre-event type)`  $\longrightarrow$  *event or #f* procedure

Same as *next-sigevent*, but returns *#f* if no appropriate event exists.

`(next-sigevent-set/no-wait set pre-event)`  $\longrightarrow$  *event or #f* procedure

Same as *next-sigevent-set*, but returns *#f* if no appropriate event exists.

As a small example, consider this piece of code that toggles the variable state by USR1 and USR2:

```
(define state #t)

(let lp ((sigevent (most-recent-sigevent)))
  (let ((next (next-sigevent sigevent interrupt/usr1)))
    (set! state #f)
    (let ((next (next-sigevent next interrupt/usr2)))
      (set! state #t)
      (lp next)))))
```

## 9.5 Interaction between threads and process state

In Unix, a number of resources are global to the process: signal handlers, working directory, *umask*, environment, user and group ids. Modular programming is difficult in the context of this global state and for concurrent programming things get even worse. Section 9.4 presents how *scsh* turns the global,

asynchronous signals handlers into modular, synchronous sigevents. Concurrent programming also benefit from sigevents as every thread may chase down the sigevent chain separately.

Scsh treats the working directory, umask, environment, and the effective user/group ID as thread-local resources. The initial value of the resources is determined by the way a thread is started: `spawn` assigns the initial values whereas `fork-thread` adopts the values of its parent. Here is a detailed description of the whole facility:

- The procedures to access and modify the resources remain as described in the previous chapters (`cwd` and `chdir`, `umask` and `set-umask`, `getenv` and `putenv`).
- Every thread receives its own copy of each resource.
- If `spawn` is used to start a new thread, the values of the resources are the same as they where at the start of `scsh`.
- The procedure

(`fork-thread thunk`)  $\longrightarrow$  *undefined* procedure  
 from the structure `thread-fluids` starts a thread which inherits the values of all resources from its parent. This behaviour is similar to what happens at process forking.

- The actual process state is updated only when necessary, i.e. on access or modification but not on context switch from one thread to another.

(`spoon thunk`)  $\longrightarrow$  *undefined* procedure  
 This is just an alias for `fork-thread` suggested by Alan Bawden.

For user and group identities arbitrary changing is not possible. Therefore they remain global process state: If a thread changes one of these values, all other threads see the new value. Consequently, `scsh` does not provide `with-uid` and friends.



## Chapter 10

# Miscellaneous routines

### 10.1 Integer bitwise ops

<code>(arithmetic-shift <i>i j</i>)</code>	$\longrightarrow$ <i>integer</i>	procedure
<code>(bitwise-and <i>i j</i>)</code>	$\longrightarrow$ <i>integer</i>	procedure
<code>(bitwise-ior <i>i j</i>)</code>	$\longrightarrow$ <i>integer</i>	procedure
<code>(bitwise-not <i>i</i>)</code>	$\longrightarrow$ <i>integer</i>	procedure
<code>(bitwise-xor <i>i j</i>)</code>	$\longrightarrow$ <i>integer</i>	procedure

These operations operate on integers representing semi-infinite bit strings, using a 2's-complement encoding.

`arithmetic-shift` shifts *i* by *j* bits. A left shift is *j* > 0; a right shift is *j* < 0.

### 10.2 Password encryption

<code>(crypt <i>key salt</i>)</code>	$\longrightarrow$ <i>encrypted value</i>	procedure
--------------------------------------	--	-----------

Decrypts *key* by directly calling the `crypt` function using *salt* to perturb the hashing algorithm. *Salt* must be a two-character string consisting of digits, alphabetic characters, "." or "\". The length of *key* may be at most eight.

### 10.3 Dot-Locking

Section 3.2.8 already points out that POSIX's file locks are almost useless in practice. To bypass this restriction other advisory locking mechanisms, based only on standard file operations, were invented. One of them is the so-called

*dot-locking* scheme where the lock of *file-name* is represented by the file *file-name.lock*. Care is taken that only one process may generate the lock for a given file.

Here is *scsh*'s interface to dot-locking:

(*obtain-dot-lock file-name [interval retry-number stale-time]*)  $\longrightarrow$  *boolean* procedure

Tries to obtain the lock for *file-name*. If the file is already locked, the thread sleeps for *interval* seconds (default is 1) before it retries. If the lock cannot be obtained after *retry-number* attempts, the procedure returns *#f*, otherwise *#t*. The default value of *retry-number* is *#f* which corresponds to an infinite number of retries.

If *stale-time* is non-*#f*, it specifies the minimum age a lock may have (in seconds) before it is considered *stale*. *Obtain-dot-lock* attempts to delete stale locks. If it was successful obtaining a lock after breaking it, *obtain-dot-lock* returns broken. If *stale-time* is *#f*, *obtain-dot-lock* never considers a lock stale. The default for *stale-time* is 300.

Note that it is possible that *obtain-dot-lock* breaks a lock but nevertheless fails to obtain it otherwise. If it is necessary to handle this case specially, use *break-dot-lock* directly (see below) rather than specifying a non-*#f* *stale-time*

(*break-dot-lock file-name*)  $\longrightarrow$  *undefined* procedure

Breaks the lock for *file-name* if one exists. Note that breaking a lock does *not* imply a subsequent *obtain-dot-lock* will succeed, as another party may have acquired the lock between *break-dot-lock* and *obtain-dot-lock*.

(*release-dot-lock file-name*)  $\longrightarrow$  *boolean* procedure

Releases the lock for *file-name*. On success, *release-dot-lock* returns *#t*, otherwise *#f*. Note that this procedure can also be used to break the lock for *file-name*.

(*with-dot-lock\* file-name thunk*)  $\longrightarrow$  *value(s) of thunk* procedure

(*with-dot-lock file-name body ...*)  $\longrightarrow$  *value(s) of body* syntax

The procedure *with-dot-lock\** obtains the requested lock, and then calls (*thunk*). When *thunk* returns, the lock is released. A non-local exit (*e.g.*, throwing to a saved continuation or raising an exception) also causes the lock to be released.

After a normal return from *thunk*, its return values are returned by *with-dot-lock\**. The *with-dot-lock* special form is equivalent syntactic sugar.

## 10.4 Syslog facility

(Note: the functionality presented in this section is still somewhat experimental and thus subject to interface changes.)

The procedures in this section provide access to the 4.2BSD syslog facility present in most POSIX systems. The functionality is in a structure called `syslog`. There's an additional structure `syslog-channels` documented below. The `scsh` interface to the syslog facility differs significantly from that of the Unix library functionality in order to support multiple simultaneous connections to the syslog facility.

Log messages carry a variety of parameters beside the text of the message itself, namely a set of options controlling the output format and destination, the facility identifying the class of programs the message is coming from, an identifier specifying the concrete program, and the level identifying the importance of the message. Moreover, a log mask can prevent messages at certain levels to be actually sent to the syslog daemon.

### Log options

A log option specifies details of the I/O behavior of the syslog facility. A syslog option is an element of a finite type (see the Scheme 48 manual) constructed by the `syslog-option` macro. The syslog facility works with sets of options which are represented as enum sets (see the Scheme 48 manual).

<code>(syslog-option <i>option-name</i>)</code>	$\longrightarrow$ <i>option</i>	syntax
<code>(syslog-option? <i>x</i>)</code>	$\longrightarrow$ <i>boolean</i>	procedure
<code>(make-syslog-options <i>list</i>)</code>	$\longrightarrow$ <i>options</i>	procedure
<code>(syslog-options <i>option-name</i> ...)</code>	$\longrightarrow$ <i>options</i>	syntax
<code>(syslog-options? <i>x</i>)</code>	$\longrightarrow$ <i>boolean</i>	procedure

`Syslog-option` constructs a log option from the name of an option. (The possible names are listed below.) `Syslog-option?` is a predicate for log options. Options are comparable using `eq?`. `Make-syslog-options` constructs a set of options from a list of options. `Syslog-options` is a macro which expands into an expression returning a set of options from names. `Syslog-options?` is a predicate for sets of options.

Here is a list of possible names of syslog options:

`console` If syslog cannot pass the message to syslogd it will attempt to write the message to the console.

**delay** Delay opening the connection to syslogd immediately until the first message is logged.

**no-delay** Open the connection to syslogd immediately. Normally the open is delayed until the first message is logged. Useful for programs that need to manage the order in which file descriptors are allocated.

**NOTA BENE:** The **delay** and **no-delay** options are included for completeness, but do not have the expected effect in the present Scheme interface: Because the Scheme interface has to multiplex multiple simultaneous connections to the syslog facility over a single one, open and close operations on that facility happen at unpredictable times.

**log-pid** Log the process id with each message: useful for identifying instantiations of daemons.

### Log facilities

A log facility identifies the originator of a log message from a finite set known to the system. Each originator is identified by a name:

(syslog-facility *facility-name*)  $\longrightarrow$  *facility* syntax

(syslog-facility? *x*)  $\longrightarrow$  *boolean* procedure

Syslog-facility is macro that expands into an expression returning a facility for a given name. Syslog-facility? is a predicate for facilities. Facilities are comparable via eq?.

Here is a list of possible names of syslog facilities:

**authorization** The authorization system: login, su, getty, etc.

**cron** The cron daemon.

**daemon** System daemons, such as routed, that are not provided for explicitly by other facilities.

**kernel** Messages generated by the kernel.

**lpr** The line printer spooling system: lpr, lpc, lpd, etc.

**mail** The mail system.

**news** The network news system.

**user** Messages generated by random user processes.

**uucp** The uucp system.

**local0 local1 local2 local3 local4 local5 local6 local7** Reserved for local use.

## Log levels

A log level identifies the importance of a message from a fixed set of possible levels.

(syslog-level *level-name*)  $\longrightarrow$  *level* syntax

(syslog-level? *x*)  $\longrightarrow$  *boolean* procedure

Syslog-level is macro that expands into an expression returning a facility for a given name. Syslog-level? is a predicate for facilities. Levels are comparable via eq?.

Here is a list of possible names of syslog levels:

**emergency** A panic condition. This is normally broadcast to all users.

**alert** A condition that should be corrected immediately, such as a corrupted system database.

**critical** Critical conditions, e.g., hard device errors.

**error** Errors.

**warning** Warning messages.

**notice** Conditions that are not error conditions, but should possibly be handled specially.

**info** Informational messages.

**debug** Messages that contain information normally of use only when debugging a program.

## Log masks

A log masks can mask out log messages at a set of levels. A log mask is an enum set of log levels.

(make-syslog-mask *list*)  $\longrightarrow$  *mask* procedure

(syslog-mask *level-name* ...)  $\longrightarrow$  *mask* syntax

syslog-mask-all mask

(syslog-mask-upto *level*)  $\longrightarrow$  *mask* procedure

(syslog-mask? *x*)  $\longrightarrow$  *boolean* procedure

Make-syslog-mask constructs a mask from a list of levels. Syslog-mask is a macro which constructs a mask from names of levels. Syslog-mask-all is a predefined log mask containing all levels. Syslog-mask-upto returns a mask consisting of all levels up to and including a certain level, starting with emergency.

## Logging

Scheme 48 dynamically maintains implicit connections to the syslog facility specifying a current identifier, current options, a current facility and a current log mask. This implicit connection is held in a thread fluid (see Section 9.5). Hence, every thread maintains its own implicit connection to syslog. Note that the connection is not implicitly preserved across a `spawn`, but it is preserved across a `fork-thread`:

(with-syslog-destination *string options facility mask thunk*)  $\longrightarrow$  *value* procedure

(set-syslog-destination! *string options facility mask*)  $\longrightarrow$  *undefined* procedure

`With-syslog-destination` dynamically binds parameters of the implicit connection to the syslog facility and runs *thunk* within those parameter bindings, returning what *thunk* returns. Each of the parameters may be `#f` in which case the previous values will be used. `Set-syslog-destination!` sets the parameters of the implicit connection of the current thread.

(syslog *level message*)  $\longrightarrow$  *undefined* procedure

(syslog *level message [string options syslog-facility]*)  $\longrightarrow$  *undefined* procedure

`Syslog` actually logs a message. Each of the parameters of the implicit connection (except for the log mask) can be explicitly specified as well for the current call to `syslog`, overriding the parameters of the channel. The parameters revert to their original values after the call.

## Syslog channels

The `syslog-channels` structure allows direct manipulation of syslog channels, the objects that represent connections to the syslog facility. Note that it is not necessary to explicitly open a syslog channel to do logging.

(open-syslog-channel *string options facility mask*)  $\longrightarrow$  *channel* procedure

(close-syslog-channel *channel*)  $\longrightarrow$  *undefined* procedure

(syslog *level message channel*)  $\longrightarrow$  *undefined* procedure

`Open-syslog-channel` and `close-syslog-channel` create and destroy a connection to the syslog facility, respectively. The specified form of calling `syslog` logs to the specified channel.

## 10.5 MD5 interface

Scsh provides a direct interface to the MD5 functions to compute the “fingerprint” or “message digest” of a file or string. It uses the C library written by Colin Plum.

(md5-digest-for-string *string*)  $\longrightarrow$  *md5-digest* procedure

Calculates the MD5 digest for the given string.

(md5-digest-for-port *port* [*buffer-size*])  $\longrightarrow$  *md5-digest* procedure

Reads the contents of the port and calculates the MD5 digest for it. The optional argument *buffer-size* determines the size of the port’s input buffer in bytes. It defaults to 1024 bytes.

(md5-digest? *thing*)  $\longrightarrow$  *boolean* procedure

The type predicate for MD5 digests: *md5-digest?* returns true if and only if *thing* is a MD5 digest.

(md5-digest->number *md5-digest*)  $\longrightarrow$  *number* procedure

Returns the number corresponding to the MD5 digest.

(number->md5-digest *number*)  $\longrightarrow$  *md5-digest* procedure

Creates a MD5 digest from a number.

(make-md5-context)  $\longrightarrow$  *md5-context* procedure

(md5-context? *thing*)  $\longrightarrow$  *boolean* procedure

(update-md5-context! *md5-context* *string*)  $\longrightarrow$  *undefined* procedure

(md5-context->md5-digest *md5-context*)  $\longrightarrow$  *md5-digest* procedure

These procedures provide a low-level interface to the library. A *md5-context* stores the state of a MD5 computation, it is created by *make-md5-context*, its type predicate is *md5-context?*. The procedure *update-md5-context!* extends the *md5-context* by the given string. Finally, *md5-context->md5-digest* returns the *md5-digest* for the *md5-context*. With these procedures it is possible to incrementally add strings to a *md5-context* before computing the digest.

## Chapter 11

# Running scsh

Scsh is currently implemented on top of Scheme 48, a freely-available Scheme implementation written by Jonathan Rees and Richard Kelsey. Scheme 48 uses a byte-code interpreter for good code density, portability and medium efficiency. It is R5RS. It also has a module system designed by Jonathan Rees.

Scsh's design is not Scheme 48 specific, although the current implementation is necessarily so. Scsh is intended to be implementable in other Scheme implementations. The Scheme 48 virtual machine that scsh uses is a specially modified version; standard Scheme 48 virtual machines cannot be used with the scsh heap image.

There are several different ways to invoke scsh. You can run it as an interactive Scheme system, with a standard read-eval-print interaction loop. Scsh can also be invoked as the interpreter for a shell script by putting a `#!/usr/local/bin/scsh -s` line at the top of the shell script.

Descending a level, it is also possible to invoke the underlying virtual machine byte-code interpreter directly on dumped heap images. Scsh programs can be pre-compiled to byte-codes and dumped as raw, binary heap images. Writing heap images strips out unused portions of the scsh runtime (such as the compiler, the debugger, and other complex subsystems), reducing memory demands and saving loading and compilation times. The heap image format allows for an initial `#!/usr/local/lib/scsh/scshvm` trigger on the first line of the image, making heap images directly executable as another kind of shell script.

Finally, scsh's static linker system allows dumped heap images to be compiled to a raw Unix `a.out(5)` format, which can be linked into the text section of the vm binary. This produces a true Unix executable binary file. Since the byte codes comprising the program are in the file's text section, they are not traced or copied by the garbage collector, do not occupy space in the vm's heap, and



do not need to be loaded and linked at startup time. This reduces the program's startup time, memory requirements, and paging overhead.

This chapter will cover these various ways of invoking scsh programs.

## 11.1 Scsh command-line switches

When the scsh top-level starts up, it scans the command line for switches that control its behaviour. These arguments are removed from the command line; the remaining arguments can be accessed as the value of the scsh variable `command-line-arguments`.

### 11.1.1 Scripts and programs

The scsh command-line switches provide sophisticated support for the authors of shell scripts and programs; they also allow the programmer to write programs that use the Scheme 48 module system.

There is a difference between a *script*, which performs its action *as it is loaded*, and a *program*, which is loaded/linked, and then performs its action by having control transferred to an entry point (*e.g.*, the `main()` function in C programs) that was defined by the load/link operation.

A *script*, by the above definition, cannot be compiled by the simple mechanism of loading it into a scsh process and dumping out a heap image—it executes as it loads. It does not have a top-level `main()`-type entry point.

It is more flexible and useful to implement a system as a program than as a script. Programs can be compiled straightforwardly; they can also export procedural interfaces for use by other Scheme packages. However, scsh supports both the script and the program style of programming.

### 11.1.2 Inserting interpreter triggers into scsh programs

When Unix tries to execute an executable file whose first 16 bits are the character pair “#!”, it treats the file not as machine-code to be directly executed by the native processor, but as source code to be executed by some interpreter. The interpreter to use is specified immediately after the “#!” sequence on the first line of the source file (along with one optional initial argument). The kernel reads in the name of the interpreter, and executes that instead. The interpreter is passed the source filename as its first argument, with the original arguments following. Consult the Unix man page for the `exec` system call for more information.

Scsh allows Scheme programs to have these triggers placed on their first line. Scsh treats the character sequence “#!” as a block-comment sequence,<sup>1</sup> and skips all following characters until it reads the comment-terminating sequence newline/exclamation-point/sharp-sign/newline (*i.e.*, the sequence “!#” occurring on its own line).

In this way, the programmer can arrange for an initial

```
#!/usr/local/bin/scsh -s
!#
```

header appearing in a Scheme program to be ignored when the program is loaded into scsh.

### 11.1.3 Module system

Scsh uses the Scheme 48 module system, which defines *packages*, *structures*, and *interfaces*.

**Package** A package is an environment—that is, a set of variable/value bindings. You can evaluate Scheme forms inside a package, or load a file into a package. Packages export sets of bindings; these sets are called *structures*.

**Structure** A structure is a named view on a package—a set of bindings. Other packages can *open* the structure, importing its bindings into their environment. Packages can provide more than one structure, revealing different portions of the package’s environment.

**Interface** An interface is the “type” of a structure. An interface is the set of names exported by a structure. These names can also be marked with other static information (*e.g.*, advisory type declarations, or syntax information).

More information on the the Scheme 48 module system can be found in the file `module.ps` in the `doc` directory of the Scheme 48 and scsh releases.

Programming Scheme with a module system is different from programming in older Scheme implementations, and the associated development problems are consequently different. In Schemes that lack modular abstraction mechanisms, everything is accessible; the major problem is preventing name-space conflicts. In Scheme 48, name-space conflicts vanish; the major problem is that not all bindings are accessible from every place. It takes a little extra work to specify what packages export which values.

---

<sup>1</sup>Why a block-comment instead of an end-of-line delimited comment? See the section on meta-args.

It may take you a little while to get used to the new style of program development. Although `scsh` can be used without referring to the module system at all, we recommend taking the time to learn and use it. The effort will pay off in the construction of modular, factorable programs.

### Module warning

Most `scsh` programs will need to import from the `scheme` structure as well as from the `scsh` structure. However, putting both of these structures in the same open clause is a bad idea because the structures `scheme` and `scsh` export some names of I/O functions in common but with different definitions. The current implementation of the module system does not recognize this as an error but silently overwrites the exports of one structure with the exports of the other. If the `scheme` structure overwrites the exports of the `scsh` structures the program will access the R<sup>5</sup>RS definitions of the I/O functions which is not what you want.

Previous versions of this manual suggested to list `scheme` and `scsh` in a specific order in the open clause of a structure to ensure that the definitions from `scsh` overwrite the ones from `scheme`. This approach is error-prone and fragile: A simple change in the implementation of the module system will render thousands of programs useless. Starting with release 0.6.3 `scsh` provides a better means to deal with this problem: the structure `scheme-with-scsh` provides all the exports of the modules `scheme` and `scsh` but exports the right denotations for the I/O functions in question. To make a long story short:

Scsh programs should open the structure `scheme-with-scsh` if they need access to the exports of `scheme` and `scsh`.

For programs which should run in versions of `scsh` prior to release 0.6.3, programmers should make sure to always put the `scsh` reference first.

#### 11.1.4 Switches

The `scsh` top-level takes command-line switches in the following format:

```
scsh [meta-arg] [switchi ...] [end-option arg1 ... argn]
```

where

<i>meta-arg:</i>	<i>\ script-file-name</i>	
<i>switch:</i>	<i>-e entry-point</i>	Specify top-level entry-point.
	<i>-o structure</i>	Open structure in current package.
	<i>-m structure</i>	Switch to package.
	<i>-n new-package</i>	Switch to new package.
	<i>-lm module-file-name</i>	Load module into config package.
	<i>-le exec-file-name</i>	Load module into exec package.
	<i>-l file-name</i>	Load file into current package.
	<i>-ll module-file-name</i>	As in <i>-lm</i> , but search the library path list.
	<i>-lel exec-file-name</i>	As in <i>-le</i> , but search the library path list.
	<i>+lp dir</i>	Add <i>dir</i> to front of library path list.
	<i>lp+ dir</i>	Add <i>dir</i> to end of library path list.
	<i>+lpe dir</i>	<i>+lp</i> , with <i>env var</i> and <i>~user</i> expansion.
	<i>lpe+ dir</i>	<i>lp+</i> , with <i>env var</i> and <i>~user</i> expansion.
	<i>+lpsd</i>	Add script-file's <i>dir</i> to front of path list.
	<i>lpsd+</i>	Add script-file's <i>dir</i> to end of path list.
	<i>-lp-clear</i>	Clear library path list to ().
	<i>-lp-default</i>	Reset library path list to system default.
	<i>-ds</i>	Do script.
	<i>-dm</i>	Do script module.
	<i>-de</i>	Do script exec.
<i>end-option:</i>	<i>-s script</i>	
	<i>-sfd num</i>	
	<i>-c exp</i>	
	<i>--</i>	

These command-line switches essentially provide a little linker language for linking a shell script or a program together with Scheme 48 modules or Scheme 48 exec programs<sup>2</sup>. The command-line processor serially opens structures and loads code into a given package. Switches that side-effect a package operate on a particular “current” package; there are switches to change this package. (These switches provide functionality equivalent to the interactive *,open*, *,load*, *,in* and *,new* commands.) Except where indicated, switches specify actions that are executed in a left-to-right order. The initial current package is the user package, which is completely empty and opens (imports the bindings of) the R5RS and scsh structures.

If the Scheme process is started up in an interactive mode, then the current package in force at the end of switch scanning is the one inside which the

<sup>2</sup>See the Section “Command programs” in the Scheme 48 manual for a description of the exec language.

interactive read-eval-print loop is started.

The command-line switch processor works in two passes: it first parses the switches, building a list of actions to perform, then the actions are performed serially. The switch list is terminated by one of the *end-option* switches. The  $arg_i$  arguments occurring after an end-option switch are passed to the scsh program as the value of `command-line-arguments` and the tail of the list returned by `(command-line)`. That is, an *end-option* switch separates switches that control the scsh “machine” from the actual arguments being passed to the scsh program that runs on that machine.

The following switches and end options are defined:

- `-o struct`  
Open the structure in the current package.
- `-n package`  
Make and enter a new package. The package has an associated structure named *package* with an empty export list. If *package* is the string “#f”, the new package is anonymous, with no associated named structure.  
  
The new package initially opens no other structures, not even the R5RS bindings. You must follow a “`-n foo`” switch with “`-o scheme`” to access the standard identifiers such as `car` and `define`.
- `-m struct`  
Change the current package to the package underlying structure *struct*. (The `-m` stands for “module.”)
- `-lm module-file-name`  
Load the specified file into scsh’s config package — the file must contain source written in the Scheme 48 module language (“load module”). Does not alter the current package.
- `-le exec-file-name`  
Load the specified file into scsh’s exec package — the file must contain source written in the Scheme 48 exec language (“load exec”). Does not alter the current package.
- `-l file-name`  
Load the specified file into the current package.
- `-c exp`  
Evaluate expression *exp* in the current package and exit. This is called `-c` after a common shell convention (see `sh` and `csh`). The expression is evaluated in the the current package (and hence is affected by `-m`’s and `-n`’s.)

When the `scsh` top-level constructs the `scsh` command-line in this case, it takes "`scsh`" to be the program name. This switch terminates argument scanning; following args become the tail of the command-line list.

- `-e entry-point`

Specify an entry point for a program. The *entry-point* is a variable that is taken from the current package in force at the end of switch evaluation. The entry point does not have to be exported by the package in a structure; it can be internal to the package. The top level passes control to the entry point by applying it to the command-line list (so programs executing in private packages can reference their command-line arguments without opening the `scsh` package to access the `(command-line)` procedure). Note that, like the list returned by the `(command-line)` procedure, the list passed to the entry point includes the name of the program being executed (as the first element of the list), not just the arguments to the program.

A `-e` switch can occur anywhere in the switch list, but it is the *last* action performed by switch scanning if it occurs. (We violate ordering here as the shell-script `#!` mechanism prevents you from putting the `-e` switch last, where it belongs.)

- `-s script`

Specify a file to load. A `-ds` (do-script), `-dm` (do-module), or `-de` (do-exec) switch occurring earlier in the switch list gives the place where the script should be loaded. If there is no `-ds`, `-dm`, or `-de` switch, then the script is loaded at the end of switch scanning, into the module that is current at the end of switch scanning.

We use the `-ds` switch to violate left-to-right switch execution order as the `-s` switch is *required* to be last (because of the `#!` machinery), independent of when/where in the switch-processing order it should be loaded.

When the `scsh` top-level constructs the `scsh` command-line in this case, it takes *script* to be the program name. This switch terminates switch parsing; following args are ignored by the switch-scanner and are passed through to the program as the tail of the command-line list.

- `-sfd num`

Loads the script from file descriptor *num*. This switch is like the `-s` switch, except that the script is loaded from one of the process' open input file descriptors. For example, to have the script loaded from standard input, specify `-sfd 0`.

- `--`

Terminate argument scanning and start up `scsh` in interactive mode. If the argument list just runs out, without either a terminating `-s` or

-- arg, then scsh also starts up in interactive mode, with an empty `command-line-arguments` list (for example, simply entering `scsh` at a shell prompt with no args at all).

When the `scsh` top-level constructs the `scsh` command-line in this case, it takes "`scsh`" to be the program name. This switch terminates switch parsing; following args are ignored by the switch-scanner and are passed through to the program as the tail of the command-line list.

- `-ds`  
Specify when to load the script ("`do-script`"). If this switch occurs, the switch list *must* be terminated by a `-s script` switch. The script is loaded into the package that is current at the `-ds` switch.

- `-dm`  
As above, but the current module is ignored. The script is loaded into the `config` package ("`do-module`"), and hence must be written in the Scheme 48 module language. This switch doesn't affect the current module—after executing this switch, the current module is the same as as it was before.

This switch is provided to make it easy to write shell scripts in the Scheme 48 module language.

- `-de`  
As above, but the current module is ignored. The script is loaded into the `exec` package ("`do-exec`"), and hence must be written in the Scheme 48 `exec` language.

This switch is provided to make it easy to write shell scripts in the Scheme 48 `exec` language.

- `-ll module-file-name`  
Load library module into `config` package. This is just like the `-lm` switch, except that it searches the library-directory path list for the file to load.

Specifically, it means: search through the *library-directories* list of directories looking for a module file of the given name, and load it in.

The *library-directories* list defaults to ("`$prefix/lib/scsh/modules/`"). Starting with version 0.6.5 the option `--with-lib-dirs-list` of the `configure` script changes this default value.

If the environment variable `$SCSH_LIB_DIRS` is set, it is used to determine the library search path. The value of this environment variable is treated as a sequence of s-expressions, which are "read" from the string:

- A string is treated as a directory,
- `#f` is replaced with the default list of directories.

A `$SCSH_LIB_DIRS` assignment of this form

```
SCSH_LIB_DIRS='." "/usr/contrib/lib/scsh/" #f "/home/shivers/lib/scsh"'
```

would produce this list of strings for the *library-directories* list:

```
(". " "/usr/contrib/lib/scsh/"  
"/usr/local/lib/scsh/modules/"  
"/home/shivers/lib/scsh")
```

When searching for a directory containing a given library module, nonexistent or read-protected directories are silently ignored; it is not an error to have them in the *library-directories* list.

It is a startup error if reading the `$SCSH_LIB_DIRS` environment variable causes a read error, or produces a value that isn't a list of strings or `#f`.

Directory search can be recursive. A directory name that ends with a slash is recursively searched.

- `-le1 exec-file-name`  
As above, but load the specified file into `scsh`'s `exec` package. This is just like the `-le` switch, except that it searches the *library-directories* path list for the file to load.
- `+lp lib-dir,lp+ lib-dir`  
Add directory *lib-dir* to the beginning or end of the *library-directories* path list, respectively.  
*lib-dir* is a single directory. It is not split at colons or otherwise processed.
- `+lpe, lpe+`  
As above, except that `~`home-directory syntax and environment variables are expanded out.
- `-lp-clear, -lp-default`  
Set the *library-directories* path list to the empty list and the system default, respectively.

These two switches are useful if you would like to protect your script from influence by the `$SCSH_LIB_DIRS` variable.

In these cases, the `$SCSH_LIB_DIRS` environment variable is never even parsed, so a bogus value will not affect the script's execution at all.



### 11.1.5 The meta argument

The scsh switch parser takes a special command-line switch, a single backslash called the “meta-argument,” which is useful for shell scripts. If the initial command-line argument is a “\” argument, followed by a filename argument *fname*, scsh will open the file *fname* and read more arguments from the second line of this file. This list of arguments will then replace the “\” argument—*i.e.*, the new arguments are inserted in front of *fname*, and the argument parser resumes argument scanning. This is used to overcome a limitation of the #! feature: the #! line can only specify a single argument after the interpreter. For example, we might hope the following scsh script, *ekko*, would implement a simple-minded version of the Unix *echo* program:

```
#!/usr/local/bin/scsh -e main -s
!#
(define (main args)
  (map (λ (arg) (display arg) (display " "))
       (cdr args))
  (newline))
```

The idea would be that the command

```
ekko Hi there.
```

would be expanded by the `exec(2)` kernel call into

```
/usr/local/bin/scsh -e main -s ekko Hi there.
```

In theory, this would cause scsh to start up, load in file *ekko*, call the entry point on the command-line list

```
(main '("ekko" "Hi" "there."))
```

and exit.

Unfortunately, the Unix `exec(2)` syscall’s support for scripts is not very general or well-designed. It will not handle multiple arguments; the #! line is usually required to contain no more than 32 characters; it is not recursive. If these restrictions are violated, most Unix systems will not provide accurate error reporting, but either fail silently, or simply incorrectly implement the desired functionality. These are the facts of Unix life.

In the *ekko* example above, our #! trigger line has three arguments (“-e”, “main”, and “-s”), so it will not work. The meta-argument is how we work around this problem. We must instead invoke the scsh interpreter with the single \ argument, and put the rest of the arguments on line two of the program. Here’s the correct program:

```
#!/usr/local/bin/scsh \
-e main -s
!#
(define (main args)
  (map (λ (arg) (display arg) (display " "))
       (cdr args))
  (newline))
```

Now, the invocation starts as

```
ekko Hi there.
```

and is expanded by `exec(2)` into

```
/usr/local/bin/scsh \ ekko Hi there.
```

When `scsh` starts up, it expands the “\” argument into the arguments read from line two of `ekko`, producing this argument list:

```
-e main -s ekko Hi there.
      ↑
Expanded from \ ekko
```

With this argument list, processing proceeds as we intended.

## Secondary argument syntax

`Scsh` uses a very simple grammar to encode the extra arguments on the second line of the `scsh` script. The only special characters are space, tab, newline, and backslash.

- Each space character terminates an argument. This means that two spaces in a row introduce an empty-string argument.
- The tab character is not permitted (unless you quote it with the backslash character described below). This is to prevent the insidious bug where you believe you have six space characters, but you really have a tab character, and *vice-versa*.
- The newline character terminates an argument, like the space character, and also terminates the argument sequence. This means that an empty line parses to the singleton list whose one element is the empty string: `("`). The grammar doesn’t admit the empty list.

- The backslash character is the escape character. It escapes backslash, space, tab, and newline, turning off their special functions, and allowing them to be included in arguments. The ANSI C escape sequences (`\b`, `\n`, `\r` and `\t`) are also supported; these also produce argument-constituents—`\n` doesn't act like a terminating newline. The escape sequence `\nnn` for *exactly* three octal digits reads as the character whose ASCII code is *nnn*. It is an error if backslash is followed by just one or two octal digits: `\3Q` is an error. Octal escapes are always constituent chars. Backslash followed by other chars is not allowed (so we can extend the escape-code space later if we like).

You have to construct these line-two argument lines carefully. In particular, beware of trailing spaces at the end of the line—they'll give you extra trailing empty-string arguments. Here's an example:

```
#!/bin/interpreter \
foo bar quux\ yow
```

would produce the arguments

```
("foo" "bar" "" "quux yow")
```

### 11.1.6 Examples

- `scsh -dm -m myprog -e top -s myprog.scm`  
Load `myprog.scm` into the `config` package, then shift to the `myprog` package and call `(top '("myprog.scm"))`, then exit. This sort of invocation is typically used in `#!` script lines (see below).
- `scsh -c '(display "Hello, world.")'`  
A simple program.
- `scsh -o bigscheme`  
Start up interactively in the `user` package after opening structure `bigscheme`.
- `scsh -o bigscheme -- Three args passed`  
Start up interactively in the `user` package after opening `bigscheme`. The `command-line-args` variable in the `scsh` package is bound to the list `("Three" "args" "passed")`, and the `(command-line)` procedure returns the list `("scsh" "Three" "args" "passed")`.
- Program `ekko`  
This shell script, called `ekko`, implements a version of the Unix `echo` program:

```
#!/usr/local/bin/scsh -s
!#
(for-each (λ (arg) (display arg) (display " "))
  command-line-args)
```

Note this short program is an example of a *script*—it executes as it loads. The Unix rule for executing `#!` shell scripts causes

```
ekko Hello, world.
```

to expand as

```
/usr/local/bin/scsh -s ekko Hello, world.
```

- Program `ekko`

This is the same program, *not* as a script. Writing it this way makes it possible to compile the program (and then, for instance, dump it out as a heap image).

```
#!/usr/local/bin/scsh \
-e top -s
!#
(define (top args)
  (for-each (λ (arg) (display arg) (display " "))
    (cdr args)))
```

The `exec(2)` expansion of the `#!` line together with the `scsh` expansion of the “`\ ekko`” meta-argument (see section 11.1.5) gives the following command-line expansion:

```
ekko Hello, world.
⇒ /usr/local/bin/scsh \ ekko          Hello, world.
⇒ /usr/local/bin/scsh -e top -s ekko Hello, world.
```

- Program `sort`

This is a program to replace the Unix `sort` utility—sorting lines read from `stdin`, and printing the results on `stdout`. Note that the source code defines a general sorting package, which is useful (1) as a Scheme module exporting sort procedures to other Scheme code, and (2) as a standalone program invoked from the `top` procedure.

```

#!/usr/local/bin/scsh \
-dm -m sort-toplevel -e top -s
!#

;;; This is a sorting module. TOP procedure exports
;;; the functionality as a Unix program akin to sort(1).
(define-structures ((sort-struct (export sort-list
                                       sort-vector!))
                   (sort-toplevel (export top)))
  (open scheme)

  (begin (define (sort-list elts <=) ...)
         (define (sort-vec! vec <=) ...)

         ;; Parse the command line and
         ;; sort stdin to stdout.
         (define (top args)
           ...)))

```

The expansion below shows how the command-line scanner (1) loads the config file `sort` (written in the Scheme 48 module language), (2) switches to the package underlying the `sort-toplevel` structure, (3) calls `(top '("sort" "foo" "bar"))` in the package, and finally (4) exits.

```

sort foo bar
⇒ /usr/local/bin/scsh \ sort                                foo bar
⇒ /usr/local/bin/scsh -dm -m sort-toplevel -e top -s sort foo bar

```

An alternate method would have used a

```
-n #f -o sort-toplevel
```

sequence of switches to specify a top-level package.

Note that the `sort` example can be compiled into a Unix program by loading the file into an `scsh` process, and dumping a heap with top-level `top`. Even if we don't want to export the `sort`'s functionality as a subroutine library, it is still useful to write the `sort` program with the module language. The command line design allows us to run this program as either an interpreted script (given the `#!` args in the header) or as a compiled heap image.

### 11.1.7 Process exit values

`Scsh` ignores the value produced by its top-level computation when determining its exit status code. If the top-level computation completed with no errors,

scsh dies with exit code 0. For example, a scsh process whose top-level is specified by a `-c exp` or a `-e entry` entry point ignores the value produced by evaluating *exp* and calling *entry*, respectively. If these computations terminate with no errors, the scsh process exits with an exit code of 0.

To return a specific exit status, use the `exit` procedure explicitly, *e.g.*,

```
scsh -c \
  "(exit (status:exit-val (run (| (fmt) (mail shivers))))))"
```

## 11.2 The scsh virtual machine

To run the Scheme 48 implementation of scsh, you run a specially modified copy of the Scheme 48 virtual machine with a scsh heap image. The scsh binary is actually nothing but a small cover program that invokes the byte-code interpreter on the scsh heap image for you. This allows you to simply start up an interactive scsh from a command line, as well as write shell scripts that begin with the simple trigger

```
#!/usr/local/bin/scsh -s
```

You can also directly execute the virtual machine, which takes its own set of command-line switches.. For example, this command starts the vm up with a 1Mword heap (split into two semispaces):

```
scshvm -o scshvm -h 1000000 -i scsh.image arg1 arg2 ...
```

The vm peels off initial vm arguments up to the `-i` heap image argument, which terminates vm argument parsing. The rest of the arguments are passed off to the scsh top-level. Scsh's top-level removes scsh switches, as discussed in the previous section; the rest show up as the value of `command-line-arguments`.

Directly executing the vm can be useful to specify non-standard switches, or invoke the virtual machine on special heap images, which can contain pre-compiled scsh programs with their own top-level procedures.

### 11.2.1 VM arguments

The vm takes arguments in the following form:

```
scshvm [meta-arg] [vm-options+] [end-option scheme-args]
```

where

```

meta-arg:  \ filename

vm-option: -h heap-size-in-words
           -s stack-size-in-words
           -o object-file-name

end-option: -i image-file-name
           --

```

The vm's meta-switch "*\ filename*" is handled the same as scsh's meta-switch, and serves the same purpose.

## VM options

The *-o object-file-name* switch tells the vm where to find relocation information for its foreign-function calls. Scsh will use a pre-compiled default if it is not specified. Scsh *must* have this information to run, since scsh's syscall interfaces are done with foreign-function calls.

The *-h* and *-s* options tell the vm how much space to allocate for the heap and stack. The heap size value is the total number of words allocated for the heap; this space is then split into two semi-spaces for Scheme 48's stop-and-copy collector.

## End options

End options terminate argument parsing. The *-i* switch is followed by the name of a heap image for the vm to execute. The *image-file-name* string is also taken to be the name of the program being executed by the VM; this name becomes the head of the argument list passed to the heap image's top-level entry point. The tail of the argument list is constructed from all following arguments.

The *--* switch terminates argument parsing without giving a specific heap image; the vm will start up using a default heap (whose location is compiled into the vm). All the following arguments comprise the tail of the list passed off to the heap image's top-level procedure.

Notice that you are not allowed to pass arguments to the heap image's top-level procedure (*e.g.*, scsh) without delimiting them with *-i* or *--* flags.

### 11.2.2 Stripped image

Besides the standard image *scsh.image* scsh also ships with the much smaller image *stripped-scsh.image*. This image contains the same code as the standard image but has almost all debugging information removed.

`stripped-scsh.image` is intended to be used with standalone programs where startup time and memory consumption count but debugging the scheme code is not that important. To use the image the VM has to be called directly and the path to the image must be given after the `-i` argument.

### 11.2.3 Inserting interpreter triggers into heap images

Scheme 48's heap image format allows for an informational header: when the vm loads in a heap image, it ignores all data occurring before the first control-L character (ASCII 12). This means that you can insert a `"#!"` trigger line into a heap image, making it a form of executable "shell script." Since the vm requires multiple arguments to be given on the command line, you must use the meta-switch. Here's an example heap-image header:

```
#!/usr/local/lib/scsh/scshvm \  
-o /usr/local/lib/scsh/scshvm -i  
... Your heap image goes here ...
```

### 11.2.4 Inserting a double-level trigger into Scheme programs

If you're a nerd, you may enjoy doing a double-level machine shift in the trigger line of your Scheme programs with the following magic:

```
#!/usr/local/lib/scsh/scshvm \  
-o /usr/local/lib/scsh/scshvm -i /usr/local/lib/scsh/scsh.image -s  
!#  
... Your Scheme program goes here ...
```

## 11.3 Compiling scsh programs

Scsh allows you to create a heap image with your own top-level procedure. Adding the pair of lines

```
#!/usr/local/lib/scsh/scshvm \  
-o /usr/local/lib/scsh/scshvm -i
```

to the top of the heap image will turn it into an executable Unix file.

You can create heap images with the following two procedures.

`(dump-scsh-program main fname)`  $\longrightarrow$  *undefined* procedure



This procedure writes out a scsh heap image. When the heap image is executed by the Scheme 48 vm, it will call the *main* procedure, passing it the vm's argument list. When *main* returns an integer value *i*, the vm exits with exit status *i*. The Scheme vm will parse command-line switches as described in section 11.2.1; remaining arguments form the tail of the command-line list that is passed to *main*. (The head of the list is the name of the program being executed by the vm.) Further argument parsing (as described for scsh in section 11.1.4) is not performed.

The heap image created by `dump-scsh-program` has unused code and data pruned out, so small programs compile to much smaller heap images.

`(dump-scsh fname)`  $\longrightarrow$  *undefined* procedure

This procedure writes out a heap image with the standard scsh top-level. When the image is resumed by the vm, it will parse and execute scsh command-line switches as described in section 11.1.4.

You can use this procedure to write out custom scsh heap images that have specific packages preloaded and start up in specific packages.

Unfortunately, Scheme 48 does not support separate compilation of Scheme files or Scheme modules. The only way to compile is to load source and then dump out a heap image. One occasionally hears rumours that this is being addressed by the Scheme 48 development team.

## 11.4 Standard file locations

Because the `scshvm` binary is intended to be used for writing shell scripts, it is important that the binary be installed in a standard place, so that shell scripts can dependably refer to it. The standard directory for the scsh tree should be `/usr/local/lib/scsh/`. Whenever possible, the vm should be located in

`/usr/local/lib/scsh/scshvm`

and a scsh heap image should be located in

`/usr/local/lib/scsh/scsh.image`

The top-level scsh program should be located in

`/usr/local/lib/scsh/scsh`

with a symbolic link to it from

`/usr/local/bin/scsh`

The Scheme 48 image format allows heap images to have `#!` triggers, so `scsh.image` should have a `#!` trigger of the following form:

```
#!/usr/local/lib/scsh/scshvm \  
-o /usr/local/lib/scsh/scshvm -i  
... heap image goes here ...
```

# Index

\*temp-file-template\*, 52  
->uid, 64  
->username, 64  
%exec, 55  
%exit, 55  
%fork, 56  
%fork/pipe, 56  
%fork/pipe+, 57  
%read-delimited!, 141  
&, 13  
&, 13  
&&, 20  
  
reduce-port, 16  
  
absolute-file-name, 109  
accept-connection, 99  
add-after, 80  
add-before, 80  
alist->env, 77  
alist-compress, 78  
alist-delete, 78  
alist-update, 78  
arg, 65  
arg\*, 65  
argv, 65  
arithmetic-shift, 161  
ascii->char, 110  
autoreap-policy, 59  
awk, 151  
  
become-session-leader, 86  
bind-listen-accept-loop, 95  
  
bind-prepare-listen-accept-loop,  
95  
bind-socket, 98  
bitwise-and, 161  
bitwise-ior, 161  
bitwise-not, 161  
bitwise-xor, 161  
break-dot-lock, 162  
bufpol/block, 38  
bufpol/line, 38  
bufpol/none, 38  
  
call-terminally, 55  
call-with-string-output-port, 26  
call/fdes, 30  
char->ascii, 110  
char-ascii?, 110  
char-blank?, 110  
char-digit?, 110  
char-graphic?, 110  
char-hex-digit?, 110  
char-iso-control?, 110  
char-letter+digit?, 110  
char-letter?, 110  
char-lower-case?, 110  
char-printing?, 110  
char-punctuation?, 110  
char-title-case?, 110  
char-upper-case?, 110  
char-whitespace?, 110  
chdir, 62  
clean-up-cres, 137  
close, 25

- close-after, 24
- close-directory-stream, 47
- close-socket, 95
- close-syslog-channel, 166
- command-line, 64
- command-line-arguments, 64
- connect-socket, 98
- connect-socket-no-wait, 98
- connect-socket-successful?, 98
- control-tty-file-name, 87
- copy-tty-info, 84
- cpu-ticks/sec, 64
- create-directory, 41
- create-fifo, 41
- create-hard-link, 41
- create-socket, 95
- create-socket-pair, 95
- create-symlink, 41
- create-temp-file, 51
- crypt, 161
- current-thread, 157
- cwd, 62
- date, 72
- date, 74
- date->string, 75
- define-record, 7
- delete-directory, 42
- delete-file, 42
- delete-filesys-object, 42
- directory-as-file-name, 106
- directory-files, 46
- disable-tty-char, 82
- drain-tty, 85
- dump-scsh, 185
- dump-scsh-program, 184
- dup, 31
- dup->fdes, 31
- dup->inport, 31
- dup->outport, 31
- enabled-interrupts, 70
- env->alist, 77
- errno-error, 21
- error-output-port, 24
- exec, 54
- exec-epf, 13
- exec-epf, 13
- exec-path, 54
- exec-path-list, 80
- exec-path-search, 55
- exec-path/env, 54
- exec/env, 54
- exit, 55
- expand-file-name, 109
- fdes->inport, 30
- fdes->outport, 30
- fdes-flags, 32
- fdes-status, 33
- field-reader, 148
- field-splitter, 144
- file-directory?, 44
- file-executable?, 46
- file-exists?, 46
- file-fifo?, 44
- file-group, 44
- file-info, 43
- file-info-directory?, 45
- file-info-executable?, 46
- file-info-fifo?, 45
- file-info-not-executable?, 46
- file-info-not-readable?, 46
- file-info-not-writable?, 46
- file-info-readable?, 46
- file-info-regular?, 45
- file-info-socket?, 45
- file-info-special?, 45
- file-info-symlink?, 45
- file-info-writable?, 46
- file-info:atime, 43
- file-info:ctime, 43
- file-info:device, 43
- file-info:gid, 43
- file-info:inode, 43
- file-info:mode, 43
- file-info:mtime, 43
- file-info:nlinks, 43

file-info:size, 43	get-lock-region, 40
file-info:type, 43	getenv, 77
file-info:uid, 43	glob, 47
file-inode, 44	glob-quote, 49
file-last-access, 44	group-info, 64
file-last-mod, 44	group-info:gid, 64
file-last-status-change, 44	group-info:members, 64
file-match, 49	group-info:name, 64
file-mode, 44	
file-name-absolute?, 106	home-dir, 109
file-name-as-directory, 105	home-directory, 80
file-name-directory, 106	home-file, 109
file-name-directory?, 105	host-info, 101
file-name-extension, 107	
file-name-non-directory?, 105	if-match, 133
file-name-nondirectory, 106	if-sre-form, 138
file-name-sans-extension, 108	infix-splitter, 144
file-nlinks, 44	internet-address->socket-address,
file-not-executable?, 45	97
file-not-exists?, 46	interrupt-handler, 71
file-not-readable?, 45	interrupt-set, 70
file-not-writable?, 45	interrupt/alarm, 68
file-owner, 44	interrupt/alrm, 68
file-readable?, 46	interrupt/chld, 68
file-regular?, 44	interrupt/cont, 68
file-size, 44	interrupt/hup, 68
file-socket?, 44	interrupt/info, 68
file-special?, 44	interrupt/int, 68
file-symlink?, 44	interrupt/io, 68
file-type, 44	interrupt/keyboard, 68
file-writable?, 46	interrupt/memory-shortage, 68
fill-in-date!, 75	interrupt/poll, 68
flush-all-ports, 39	interrupt/prof, 68
flush-submatches, 134	interrupt/pwr, 68
flush-tty/both, 85	interrupt/quit, 68
flush-tty/input, 85	interrupt/term, 68
flush-tty/output, 85	interrupt/tstp, 68
force-output, 38	interrupt/urg, 68
fork, 56	interrupt/usr1, 68
fork-pty-session, 87	interrupt/usr2, 68
fork-thread, 160	interrupt/vtalrm, 68
fork/pipe, 56	interrupt/winch, 68
fork/pipe+, 57	interrupt/xcpu, 68
format-date, 75	interrupt/xfsz, 68

- itimer, 67
- join-strings, 147
- let-match, 133
- listen-socket, 99
- lock-owner-uid, 158
- lock-region, 40
- lock-region/no-block, 40
- lock-region:end, 39
- lock-region:exclusive?, 39
- lock-region:len, 39
- lock-region:proc, 39
- lock-region:start, 39
- lock-region:whence, 39
- lock-region?, 39
- lock?, 157
- make-char-port-filter, 20
- make-date, 73
- make-lock, 157
- make-lock-region, 40
- make-md5-context, 167
- make-placeholder, 158
- make-pty-generator, 88
- make-re-char-set, 136
- make-re-choice, 136
- make-re-dsm, 136
- make-re-repeat, 136
- make-re-seq, 135
- make-re-string, 136
- make-re-submatch, 136
- make-regexp, 129
- make-string-input-port, 26
- make-string-output-port, 26
- make-string-port-filter, 20
- make-syslog-mask, 165
- make-syslog-options, 163
- make-tty-info, 84
- match-cond, 133
- match:end, 130
- match:start, 130
- match:substring, 130
- maybe-obtain-lock, 158
- md5-context->md5-digest, 167
- md5-context?, 167
- md5-digest->number, 167
- md5-digest-for-port, 167
- md5-digest-for-string, 167
- md5-digest?, 167
- most-recent-sigevent, 159
- move->fdes, 30
- network-info, 101
- next-sigevent, 159
- next-sigevent-set, 159
- next-sigevent-set/no-wait, 159
- next-sigevent/no-wait, 159
- nice, 63
- number->md5-digest, 167
- obtain-dot-lock, 162
- obtain-lock, 157
- open-control-tty, 86
- open-directory-stream, 47
- open-fdes, 32
- open-file, 32
- open-input-file, 32
- open-output-file, 32
- open-pty, 87
- open-syslog-channel, 166
- parent-pid, 62
- parse-file-name, 108
- parse-sre, 138
- parse-sres, 138
- path-list->file-name, 107
- pid, 62
- pid->proc, 58
- pipe, 33
- placeholder-value, 158
- placeholder?, 158
- port->fdes, 30
- port->list, 15
- port->sexp-list, 15
- port->socket, 96
- port->string, 15
- port->string-list, 15

- port-fold, 16
- port-revealed, 30
- posix-string->regexp, 135
- priority, 63
- proc, 58
- proc:pid, 58
- proc?, 58
- process-group, 62
- process-sleep, 67
- process-sleep-until, 67
- process-times, 63
- protocol-info, 101
- pty-name->tty-name, 87
  
- re-any, 137
- re-any?, 137
- re-bol, 136
- re-bol?, 136
- re-bos, 136
- re-bos?, 136
- re-char-set, 136
- re-char-set:cset, 136
- re-char-set?, 136
- re-choice, 136
- re-choice:elts, 136
- re-choice:tsm, 136
- re-choice?, 136
- re-dsm, 136
- re-dsm:body, 136
- re-dsm:post-dsm, 136
- re-dsm:pre-dsm, 136
- re-dsm:tsm, 136
- re-dsm?, 136
- re-empty, 137
- re-empty?, 137
- re-eol, 136
- re-eol?, 136
- re-eos, 136
- re-eos?, 136
- re-nonl, 137
- re-repeat:from, 136
- re-repeat:to, 136
- re-repeat:tsm, 136
- re-repeat?, 136
  
- re-seq, 135
- re-seq:elts, 135
- re-seq:tsm, 136
- re-seq?, 135
- re-string, 136
- re-string:chars, 136
- re-string?, 136
- re-submatch:post-dsm, 136
- re-submatch:pre-dsm, 136
- re-submatch:tsm, 136
- re-submatch?, 136
- re-trivial, 137
- re-trivial?, 137
- re-tsm, 137
- read-delimited, 141
- read-delimited!, 141
- read-directory-stream, 47
- read-line, 140
- read-paragraph, 141
- read-string, 33
- read-string!, 33
- read-string!/partial, 34
- read-string/partial, 34
- read-symlink, 42
- reap-zombies, 59
- receive-message, 99
- receive-message!, 99
- receive-message!/partial, 99
- receive-message/partial, 99
- record-reader, 144
- regexp->posix-string, 135
- regexp->scheme, 138
- regexp->sre, 134
- regexp-fold, 131
- regexp-fold-right, 132
- regexp-for-each, 133
- regexp-search, 129
- regexp-search?, 129
- regexp-substitute, 130
- regexp-substitute/global, 130
- regexp?, 129, 137
- release-dot-lock, 162
- release-lock, 158
- release-port-handle, 30

- relinquish-timeslice, 157
- rename-file, 42
- replace-extension, 108
- resolve-file-name, 109
- run, 13
- run, 13
- run/collecting, 17
- run/collecting\*, 17
- run/file, 14
- run/file\*, 15
- run/port, 14
- run/port\*, 15
- run/port+proc, 17
- run/port+proc\*, 17
- run/sexp, 14
- run/sexp\*, 15
- run/sexps, 14
- run/sexps\*, 15
- run/string, 14
- run/string\*, 15
- run/strings, 14
- run/strings\*, 15
- rx, 129
- seek, 31
- select , 35
- select-port-channels, 36
- select-ports, 36
- send-message, 99
- send-message/partial, 99
- send-tty-break, 85
- service-info, 101
- set-enabled-interrupts, 70
- set-fdes-flags, 32
- set-fdes-status, 33
- set-file-group, 42
- set-file-mode, 42
- set-file-owner, 42
- set-file-times, 42
- set-gid, 63
- set-interrupt-handler, 70
- set-port-buffering, 38
- set-priority, 63
- set-process-group, 62
- set-socket-option, 100
- set-tty-info/drain, 84
- set-tty-info/flush, 84
- set-tty-info/now, 84
- set-tty-process-group, 86
- set-uid, 63
- set-umask, 62
- set-user-effective-gid, 63
- set-user-effective-uid, 63
- setenv, 77
- shutdown-socket, 99
- sigevent?, 159
- signal->interrupt, 67
- signal-process, 67
- signal-process-group, 67
- signal/abrt, 69
- signal/alm, 68
- signal/bus, 69
- signal/chld, 68
- signal/cont, 68
- signal/emt, 69
- signal/fpe, 69
- signal/hup, 68
- signal/ill, 69
- signal/info, 68
- signal/int, 68
- signal/io, 68
- signal/iot, 69
- signal/kill, 69
- signal/pipe, 69
- signal/poll, 68
- signal/prof, 68
- signal/pwr, 68
- signal/quit, 68
- signal/segv, 69
- signal/stop, 69
- signal/sys, 69
- signal/term, 68
- signal/trap, 69
- signal/tstp, 68
- signal/ttin, 69
- signal/ttou, 69
- signal/urg, 68
- signal/usr1, 68



- signal/usr2, 68
- signal/vtalm, 68
- signal/winch, 68
- signal/xcpu, 68
- signal/xfsz, 68
- simplify-file-name, 109
- simplify-regexp, 134
- skip-char-set, 142
- sleep, 157
- sloppy-suffix-splitter, 144
- socket-address->internet-address, 97
- socket-address->unix-address, 97
- socket-connect, 94
- socket-local-address, 99
- socket-option, 100
- socket-remote-address, 99
- spawn, 156
- split-file-name, 107
- spoon, 160
- sre->regexp, 134
- sre-form?, 138
- start-tty-input, 85
- start-tty-output, 85
- static-regexp?, 139
- status:exit-val, 61
- status:stop-sig, 61
- status:term-sig, 61
- stdio->stdports, 25
- stdports->stdio, 25
- stop-tty-input, 85
- stop-tty-output, 85
- string-match, 129
- string-output-port-output, 26
- substitute-env-vars, 110
- suffix-splitter, 144
- suspend, 56
- sync-file, 43
- sync-file-system, 43
- syslog, 166
- syslog-facility, 164
- syslog-facility?, 164
- syslog-level, 165
- syslog-level?, 165
- syslog-mask, 165
- syslog-mask-all, 165
- syslog-mask-upto, 165
- syslog-mask?, 165
- syslog-option, 163
- syslog-option?, 163
- syslog-options, 163
- syslog-options?, 163
- system-name, 66
- tell, 32
- temp-file-channel, 53
- temp-file-iterate, 52
- terminate-current-thread, 157
- thread-name, 157
- thread-uid, 157
- thread?, 157
- ticks/sec, 73
- time, 71
- time, 74
- time+ticks, 73
- truncate-file, 43
- tty-file-name, 82
- tty-info, 84
- tty-info record type, 82
- tty-info:control-chars, 82
- tty-info:control-flags, 82
- tty-info:input-flags, 82
- tty-info:input-speed, 82
- tty-info:local-flags, 82
- tty-info:min, 82
- tty-info:output-flags, 82
- tty-info:output-speed, 82
- tty-info:time, 82
- tty-info?, 82
- tty-name->pty-name, 87
- tty-process-group, 86
- tty?, 82
- ttyc/2-stop-bits, 92
- ttyc/carrier-flow-ctl, 92
- ttyc/char-size, 92
- ttyc/char-size5, 92
- ttyc/char-size6, 92
- ttyc/char-size7, 92

ttyc/char-size8, 92	ttyl/echo-delete-line, 93
ttyc/CTS-output-flow-ctl, 92	ttyl/echo-nl, 93
ttyc/enable-parity, 92	ttyl/enable-signals, 93
ttyc/enable-read, 92	ttyl/extended, 93
ttyc/hup-on-close, 92	ttyl/flush-output, 93
ttyc/ignore-flags, 92	ttyl/hardcopy-delete, 93
ttyc/no-modem-sync, 92	ttyl/no-flush-on-interrupt, 93
ttyc/odd-parity, 92	ttyl/no-kernel-status, 93
ttyc/RTS-input-flow-ctl, 92	ttyl/reprint-unread-chars, 93
ttychar/delayed-suspend, 89	ttyl/ttou-signal, 93
ttychar/delete-char, 89	ttyl/visual-delete, 93
ttychar/delete-line, 89	ttyl/visual-delete-line, 93
ttychar/delete-word, 89	ttyout/all-delay, 91
ttychar/discard, 89	ttyout/bs-delay, 91
ttychar/eof, 89	ttyout/bs-delay0, 91
ttychar/eol, 89	ttyout/bs-delay1, 91
ttychar/eol2, 89	ttyout/cr->nl, 90
ttychar/interrupt, 89	ttyout/cr-delay, 91
ttychar/literal-next, 89	ttyout/cr-delay0, 91
ttychar/quit, 89	ttyout/cr-delay1, 91
ttychar/reprint, 89	ttyout/cr-delay2, 91
ttychar/start, 89	ttyout/cr-delay3, 91
ttychar/status, 89	ttyout/delay-w/fill-char, 90
ttychar/stop, 89	ttyout/discard-eot, 90
ttychar/suspend, 89	ttyout/enable, 90
ttyin/7bits, 90	ttyout/expand-tabs, 90
ttyin/beep-on-overflow, 90	ttyout/ff-delay, 91
ttyin/check-parity, 90	ttyout/ff-delay0, 91
ttyin/cr->nl, 90	ttyout/ff-delay1, 91
ttyin/ignore-bad-parity-chars, 90	ttyout/fill-w/del, 90
ttyin/ignore-break, 90	ttyout/nl->crnl, 90
ttyin/ignore-cr, 90	ttyout/nl-delay, 91
ttyin/input-flow-ctl, 90	ttyout/nl-delay0, 91
ttyin/interrupt-on-break, 90	ttyout/nl-delay1, 91
ttyin/lowercase, 90	ttyout/nl-does-cr, 90
ttyin/mark-parity-errors, 90	ttyout/no-col0-cr, 90
ttyin/nl->cr, 90	ttyout/tab-delay, 91
ttyin/output-flow-ctl, 90	ttyout/tab-delay0, 91
ttyin/xon-any, 90	ttyout/tab-delay1, 91
ttyl/alt-delete-word, 93	ttyout/tab-delay2, 91
ttyl/canonical, 93	ttyout/tab-delayx, 91
ttyl/case-map, 93	ttyout/uppercase, 90
ttyl/echo, 93	ttyout/vtab-delay, 91
ttyl/echo-ctl, 93	ttyout/vtab-delay0, 91

ttyout/vtab-delay1, 91  
 umask, 62  
 uname, 66  
 uncase, 134  
 uncase-char-set, 134  
 uncase-string, 134  
 unix-address->socket-address, 97  
 unlock-region, 41  
 user-effective-gid, 63  
 user-effective-uid, 63  
 user-gid, 63  
 user-info, 64  
 user-info:gid, 64  
 user-info:home-dir, 64  
 user-info:name, 64  
 user-info:shell, 64  
 user-info:uid, 64  
 user-login-name, 63  
 user-supplementary-gids, 63  
 user-uid, 63  
 wait, 60  
 wait-any, 60  
 wait-process-group, 61  
 with-current-input-port, 24  
 with-current-input-port\*, 24  
 with-current-output-port, 24  
 with-current-output-port\*, 24  
 with-cwd, 62  
 with-cwd\*, 62  
 with-dot-lock, 162  
 with-dot-lock\*, 162  
 with-enabled-interrupts, 70  
 with-enabled-interrupts\*, 70  
 with-env, 79  
 with-env\*, 78  
 with-errno-handler, 22  
 with-errno-handler\*, 22  
 with-error-output-port, 24  
 with-error-output-port\*, 24  
 with-region-lock, 41  
 with-region-lock\*, 41  
 with-stdio-ports, 26  
 with-stdio-ports\*, 26  
 with-syslog-destination, 166  
 with-total-env, 79  
 with-total-env\*, 78  
 with-umask, 62  
 with-umask\*, 62  
 with-user-effective-gid, 63  
 with-user-effective-gid\*, 63  
 with-user-effective-uid, 63  
 with-user-effective-uid\*, 63  
 write-string, 36  
 write-string/partial, 37