

`}${Unix_Tools}`

Markus Kuhn

Computer Laboratory, University of Cambridge

<http://www.cl.cam.ac.uk/teaching/1314/UnixTools/>

Michaelmas 2013 – Part 1B

1

Brief review of Unix history

- ▶ “First Edition” developed at AT&T Bell Labs during 1968–71 by Ken Thompson and Dennis Ritchie for a PDP 11
- ▶ Rewritten in C in 1973
- ▶ Sixth Edition (1975) first widely available version
- ▶ Seventh Edition in 1979, UNIX 32V for VAX
- ▶ During 1980s independent continued development at AT&T (“System V Unix”) and Berkeley University (“BSD Unix”)
- ▶ Commercial variants (Solaris, SCO, HP/UX, AIX, IRIX, ...)
- ▶ IEEE and ISO standardisation of a *Portable Operating System Interface based on Unix (POSIX)* in 1989, later also *Single Unix Specification* by X/Open, both merged in 2001

The POSIX standard is freely available online: <http://www.unix.org/> and <http://pubs.opengroup.org/onlinepubs/9699919799/>

3

Why do we teach Unix Tools?

- ▶ Second most popular OS family (after Microsoft Windows)
- ▶ Many elements of Unix have become part of common computer science folklore, terminology & tradition over the past 25 years and influenced many other systems (including DOS/Windows)
- ▶ Many Unix tools have been ported and become popular on other platforms: full Unix environment in Apple's OS X, Cygwin
- ▶ Your future project supervisors and employers are likely to expect you to be fluent under Unix as a development environment
- ▶ Good examples for high-functionality user interfaces

This short lecture course can only give you a first overview. You need to spend at least 2–3 times as many hours with e.g. MCS Linux to

- ▶ explore the tools mentioned
- ▶ solve exercises (which often involve reading documentation to understand important details skipped in the lecture)

2

A brief history of free Unix

- ▶ In 1983, Richard Stallman (MIT) initiates a free reimplement of Unix called GNU (“GNU’s Not Unix”) leading to an editor (emacs), compiler (gcc), debugger (gdb), and numerous other tools.
- ▶ In 1991, Linus Torvalds (Helsinki CS undergraduate) starts development of a free POSIX-compatible kernel, later nicknamed *Linux*, which was rapidly complemented by existing GNU tools and contributions from volunteers and industry to form a full Unix replacement.
- ▶ In 1991, Berkeley University releases a free version of BSD Unix, after removing remaining proprietary AT&T code. Volunteer projects emerge to continue its development (FreeBSD, NetBSD, OpenBSD).
- ▶ In 2000, Apple releases *Darwin*, the now open-source core components of their OS X and iOS operating systems. Volunteer projects emerge to port many Unix tools onto Darwin (Homebrew, Fink, MacPorts, GNU Darwin, etc.).

4

Free software license concepts

- ▶ **public domain:** authors waive all copyright
- ▶ **“MIT/BSD” licences:** allow you to copy, redistribute and modify the software in any way as long as
 - you respect the identity and rights of the author (preserve copyright notice and licence terms in source code and documentation)
 - you agree not sue the author over software quality (accept exclusion of liability and warranty)
- ▶ **GNU General Public Licence (GPL):** requires in addition that any derived work distributed or published
 - must be licensed under the terms of the GPL
 - must have its source code made publicly available

Numerous refinements and variants of these licences have been written. For more information on the various opensource licence types and their philosophies: <http://opensource.org/>

5

Original Unix user interfaces

The initial I/O devices were teletype terminals . . .



Photo: Bell Labs

6

VT100 terminals

. . . and later video display terminals such as the DEC VT100, all providing 80 characters-per-line fixed-width ASCII output. Their communications protocol is still used today in graphical windowing environments via “terminal emulators” (e.g., xterm, konsole).



Photo: <http://www.catb.org/esr/writings/taou/html/>

The VT100 was the first video terminal with microprocessor, and the first to implement the ANSI X3.64 (= ECMA-48) control functions. For instance, “`Esc[7m`” activates **inverse mode** and “`Esc[0m`” returns to normal, where `Esc` is the ASCII “escape” control character (27 = 0x1B).

<http://www.vt100.net/>
<http://www.ecma-international.org/publications/standards/Ecma-048.htm>
`man console_codes`

7

Unix tools design philosophy

- ▶ Compact and concise input syntax, making full use of ASCII repertoire to minimise keystrokes
- ▶ Output format should be simple and easily usable as input for other programs
- ▶ Programs can be joined together in “pipes” and “scripts” to solve more complex problems
- ▶ Each tool originally performed a simple single function
- ▶ Prefer reusing existing tools with minor extension to rewriting a new tool from scratch
- ▶ The main user-interface software (“shell”) is a normal replaceable program without special privileges
- ▶ Support for automating routine tasks

Brian W. Kernighan, Rob Pike: The Unix Programming Environment. Prentice-Hall, 1984.

8

Unix documentation

Most Unix documentation can be read from the command line.

Classic manual sections: user commands (1), system calls (2), library functions (3), devices (4), file formats (5).

- ▶ The `man` tool searches for the manual page file (→ `$MANPATH`) and activates two further tools (`nroff` text formatter and `more` text-file viewer). Add optional section number to disambiguate:

```
$ man 3 printf      # C subroutine, not command
```

Honesty in documentation: Unix manual pages traditionally include a BUGS section.

- ▶ `xman`: X11 GUI variant, offers a table of contents
- ▶ `info`: alternative GNU hypertext documentation system
Invoke with `info` from the shell or with `C-h i` from `emacs`. Use `M(enu)` key to select topic or `[Enter]` to select hyperlink under cursor, `N(ext)/P(rev)/U(p)/D(irectory)` to navigate document tree, `Emacs` search function (`Ctrl-S`), and finally `Q(uit)`.
- ▶ Check `/usr/share/doc/` and Web for further documentation.

9

Examples of Unix tools

<code>man</code> , <code>apropos</code> , <code>xman</code> , <code>info</code> help/documentation browser	<code>make</code> project builder
<code>more</code> , <code>less</code> plaintext file viewer	<code>cmp</code> , <code>diff</code> , <code>patch</code> compare files, apply patches
<code>ls</code> , <code>find</code> list/traverse directories, search	<code>rsc</code> , <code>cvs</code> , <code>svn</code> , <code>git</code> , <code>hg</code> , <code>bzr</code> revision control systems
<code>cp</code> , <code>mv</code> , <code>rm</code> , <code>touch</code> , <code>ln</code> copy, move/rename, remove, renew files, link/shortcut files	<code>adb</code> , <code>gdb</code> debuggers
<code>mkdir</code> , <code>rmdir</code> make/remove directories	<code>awk</code> , <code>perl</code> , <code>python</code> , <code>tcl</code> scripting languages
<code>cat</code> , <code>dd</code> , <code>head</code> , <code>tail</code> concatenate/split files	<code>m4</code> , <code>cpp</code> macro processors
<code>du</code> , <code>df</code> , <code>quota</code> , <code>rquota</code> examine disk space used and free	<code>sed</code> , <code>tr</code> edit streams, replace characters
<code>ps</code> , <code>top</code> , <code>free</code> , <code>uptime</code> , <code>w</code> process table and system load	<code>sort</code> , <code>grep</code> , <code>cut</code> sort/search lines of text, extract columns
<code>vi</code> , <code>emacs</code> , <code>pico</code> interactive editors	<code>nroff</code> , <code>troff</code> , <code>tex</code> , <code>latex</code> text formatters
<code>cc</code> , <code>gcc</code> C compilers	<code>mail</code> , <code>pine</code> , <code>mh</code> , <code>exmh</code> , <code>elm</code> electronic mail user agents

10

Examples of Unix tools (cont'd)

<code>telnet</code> , <code>ftp</code> , <code>rlogin</code> , <code>finger</code> , <code>talk</code> , <code>ping</code> , <code>traceroute</code> , <code>wget</code> , <code>curl</code> , <code>ssh</code> , <code>scp</code> , <code>rsync</code> , <code>hostname</code> , <code>host</code> , <code>ifconfig</code> , <code>route</code> network tools	<code>clear</code> , <code>reset</code> clear screen, reset terminal
<code>xterm</code> VT100 terminal emulator	<code>stty</code> configure terminal driver
<code>tar</code> , <code>cpio</code> , <code>compress</code> , <code>zip</code> , <code>gzip</code> , <code>bzip2</code> file packaging and compression	<code>xv</code> , <code>display</code> , <code>ghostview</code> , <code>acroread</code> graphics file viewers
<code>echo</code> , <code>cd</code> , <code>pushd</code> , <code>popd</code> , <code>exit</code> , <code>ulimit</code> , <code>time</code> , <code>history</code> builtin shell commands	<code>xfig</code> , <code>tgif</code> , <code>gimp</code> , <code>inkscape</code> graphics drawing tools
<code>fg</code> , <code>bg</code> , <code>jobs</code> , <code>kill</code> builtin shell job control	<code>*topnm</code> , <code>pnmt*</code> , <code>[cd]jpeg</code> graphics format converters
<code>date</code> , <code>xclock</code> clocks	<code>bc</code> calculator
<code>which</code> , <code>whereis</code> locate command file	<code>passwd</code> change your password
	<code>chmod</code> change file permissions
	<code>lex</code> , <code>yacc</code> , <code>flex</code> , <code>bison</code> scanner/parser generators

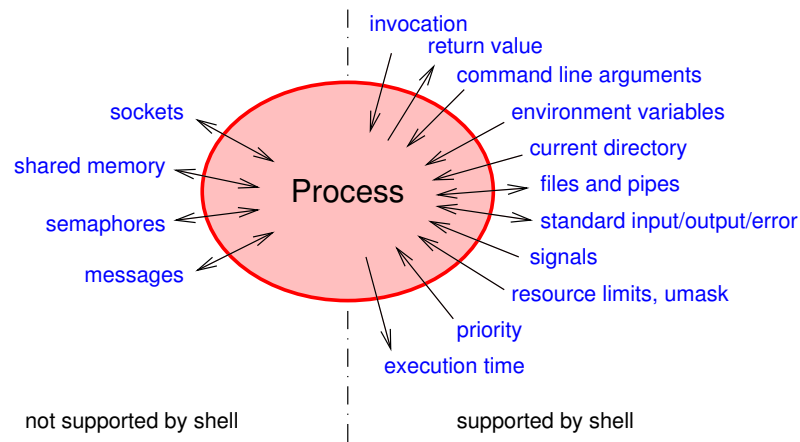
11

The Unix shell

- ▶ The user program that Unix starts automatically after a login
- ▶ Allows the user to interactively start, stop, suspend, and resume other programs and control the access of programs to the terminal
- ▶ Supports automation by executing files of commands ("shell scripts"), provides programming language constructs (variables, string expressions, conditional branches, loops, concurrency)
- ▶ Simplifies file selection via keyboard (regular expressions, file name completion)
- ▶ Simplifies entry of command arguments with editing and history functions
- ▶ Most common shell ("`sh`") developed 1975 by Stephen Bourne, modern GNU replacement is "`bash`" ("Bourne-Again SHell")

12

Unix inter-process communication mechanisms



Command-line arguments, return value, environment

A Unix C program is invoked by calling its `main()` function with:

- ▶ a list of strings `argv` as an argument
- ▶ a list of strings `environ` as a predefined global variable

```
#include <stdio.h>
extern char **environ;

int main(int argc, char **argv)
{
    int i;
    printf("Command line arguments:\n");
    for (i = 0; i < argc; i++)
        puts(argv[i]);
    printf("Environment:\n");
    for (i = 0; environ[i] != NULL; i++)
        puts(environ[i]);
    return 0;
}
```

Environment strings have the form

name=value

where *name* is free of "=".

Argument `argv[0]` is usually the name or path of the program.

Convention: `main() == 0` signals success, other values signal errors to calling process.

13

14

File descriptors

Unix processes access files in three steps:

- ▶ Provide kernel in `open()` or `creat()` system call a path name and get in return an integer "file descriptor".
- ▶ Provide in `read()`, `write()`, and `seek()` system calls an opened file descriptor along with data.
- ▶ Finally, call `close()` to release any data structures associated with an opened file (position pointer, buffers, etc.).

The `lsdf` tool lists the files currently opened by any process. Under Linux, file descriptor lists and other kernel data can be accessed via the simulated file system mounted under `/proc`.

As a convention, the shell opens three file descriptors for each process:

- ▶ 0 = standard input (for reading the data to be processed)
- ▶ 1 = standard output (for the resulting output data)
- ▶ 2 = standard error (for error messages)

Basic shell syntax – pipes

Start a program and connect the three default file descriptors `stdin`, `stdout`, and `stderr` to the terminal:

```
$ command
```

Connect `stdout` of `command1` to `stdin` of `command2` and `stdout` of `command2` to `stdin` of `command3` by forming a pipe:

```
$ command1 | command2 | command3
```

Also connects terminal to `stdin` of `command1`, to `stdout` of `command3`, and to `stderr` of all three.

Note how this function concatenation notation makes the addition of command arguments somewhat clearer compared to the mathematical notation `command3(command2(command1(arg1), arg2), arg3)`:

```
$ ls -la | sort -n -k5 | less
```

15

16

Basic shell syntax – sequences

Execute several commands (or entire pipes) in sequence:

```
$ command1 ; command2 ; command3
```

For example:

```
$ date ; host ryath.ds.cam.ac.uk
Thu Oct 31 11:48:22 GMT 2013
ryath.ds.cam.ac.uk has address 131.111.8.223
```

Conditional execution depending on success of previous command (as in logic-expression short-cut):

```
$ make ftest && ./ftest
$ ./ftest || echo 'Test failed!'
```

Return value 0 for success is interpreted as Boolean value "true", other return values for problems or failure as "false". The trivial tools `true` and `false` simply return 0 and 1, respectively.

17

File redirecting – basics

Send stdout to file

```
$ command >filename
```

Append stdout to file

```
$ command >>filename
```

Send both stdout and stderr to the same file. First redirect stdout to `filename`, then redirect stderr (file descriptor 2) to where stdout goes (target of file descriptor 1 = `&1`):

```
$ command >filename 2>&1
```

Feed stdin from file

```
$ command <filename
```

18

File redirecting – advanced

Open other file descriptors for input, output, or both

```
$ command 0<in 1>out 2>>log 3<auxin 4>auxout 5<>data
```

"Here Documents" allow us to insert data into shell scripts directly such that the shell will feed it into a command via standard input. The `<<` is followed immediately by an end-of-text marker string.

```
$ tr <<THEEND A-MN-Za-mn-z N-ZA-Mn-za-m
> Vs lbh zhfg cbfg n ehqr wbxr ba HFRARG, ebgngr gur
> nycunorg ol 13 punenpgref naq nqq n jneavat.
> THEEND
```

Redirecting to or from `/dev/tcp/hostname/port` will open a TCP socket connection:

```
{ echo "GET /~mgk25/iso-paper.c" >&3 ; cat <&3 ; } \
3<>/dev/tcp/www.c1.cam.ac.uk/80
```

The above example is a bash implementation of a simple web browser. It downloads and displays the file `http://www.c1.cam.ac.uk/~mgk25/iso-paper.c`.

Bash's `/dev/tcp/...` feature is disabled in some Linux distributions (security concerns).

19

Command-line argument conventions

Each program receives from the caller as a parameter an array of strings (`argv`). The shell places into the `argv` parameters the words entered following the command name, after several preprocessing steps have been applied first.

```
$ cp 'Lecture Timetable.pdf' lecture-timetable.pdf
$ mv *.bak old-files/
```

Command options are by convention single letters prefixed by a hyphen ("`-h`"). Unless followed by option parameters, single character flag options can often be concatenated:

```
$ ls -l -a -t
$ ls -lat
```

GNU tools offer alternatively long option names prefixed by two hyphens ("`--help`"). Arguments not starting with hyphens are typically filenames, hostnames, URLs, etc.

```
$ gcc --version
$ curl --head http://www.c1.cam.ac.uk/
```

20

Command-line argument conventions (cont'd)

The special option “--” signals in many tools that subsequent words are arguments, not options. This provides one way to access filenames starting with a hyphen:

```
$ rm -- -i
$ rm ./-i
```

The special filename “-” signals often that standard input/output should be used instead of a file.

All these are conventions that most – but not all – tools implement (usually via the getopt library function), so check the respective manual first.

The shell remains ignorant of these “-” conventions!

http://pubs.opengroup.org/onlinepubs/9699919799/basedefs/V1_chap12.html

21

Shell command-line preprocessing

A number of punctuation characters in a command line are part of the shell control syntax

```
| & ; ( ) < >
```

or can trigger special convenience substitutions before argv is handed over to the called program:

- ▶ brace expansion: {,}
- ▶ tilde expansion: ~
- ▶ parameter expansion: \$
- ▶ pathname expansion / filename matching: * ? []
- ▶ quote removal: \ ' "

22

Brace and tilde expansion

Brace expansion

Provides for convenient entry of words with repeated substrings:

```
$ echo a{b,c,d}e
abe ace ade
$ echo {mgk25,fapp2,rja14}@cam.ac.uk
mgk25@cam.ac.uk fapp2@cam.ac.uk rja14@cam.ac.uk
$ rm slides.{bak,aux,dvi,log,ps}
```

This bash extension is not required by the POSIX standard; e.g. Ubuntu Linux /bin/sh lacks it.

Tilde expansion

Provides convenient entry of home directory pathname:

```
$ echo ~pb ~/Mail/inbox
/home/pb /homes/mgk25/Mail/inbox
```

The builtin echo command simply outputs argv to stdout and is useful for demonstrating command-line expansion and for single-line text output in scripts.

23

Parameter and command expansion

Substituted with the values of shell variables

```
$ OBJFILE=skipjack.o
$ echo ${OBJFILE} ${OBJFILE%.o}.c
skipjack.o skipjack.c
$ echo ${HOME} ${PATH} ${LOGNAME}
/homes/mgk25 /bin:/usr/bin:/usr/local/bin:/usr/X11R6/bin mgk25
```

or the standard output lines of commands

```
$ which emacs
/usr/bin/emacs
$ echo $(which emacs)
/usr/bin/emacs
$ ls -l $(which emacs)
-rwxr-xr-x  2 root  system  3471896 Mar 16  2001 /usr/bin/emacs
```

Shorter alternatives: variables without braces and command substitution with grave accent (`) or, with older fonts, back quote (‘)

```
$ echo $OBJFILE
skipjack.o
$ echo `which emacs`
/usr/bin/emacs
```

24

Pathname expansion

Command-line arguments containing `?`, `*`, or `[...]` are interpreted as regular expression patterns and will be substituted with a list of all matching filenames.

- ▶ `?` stands for an arbitrary single character
- ▶ `*` stands for an arbitrary sequence of zero or more characters
- ▶ `[...]` stands for one character out of a specified set. Use `-` to specify range of characters and `^` to complement set. Certain character classes can be named within `[:...:]`.

None of the above will match a dot at the start of a filename, which is the naming convention for hidden files.

Examples:

```
*.bak [A-Za-z]*.??? [[:alpha:]]* [^A-Z].??* files/*/*.o
```

25

Job control

Start command or entire pipe as a background job, without connecting stdin to terminal:

```
$ command &
[1] 4739
$ ./testrun 2>&1 | gzip -9c >results.gz &
[2] 4741
$ ./testrun1 & ./testrun2 & ./testrun3 &
[3] 5106
[4] 5107
[5] 5108
```

Shell prints both a job number (identifying all processes in pipe) as well as process ID of last process in pipe. Shell will list all its jobs with the `jobs` command, where a `+` sign marks the last stopped (default) job.

27

Quote removal

Three quotation mechanisms are available to enter the special characters in command-line arguments without triggering the corresponding shell substitution:

- ▶ `'...'` suppresses all special character meanings
- ▶ `"..."` suppresses all special character meanings, except for `$ \ ``
- ▶ `\` suppresses all special character meanings for the immediately following character

Example:

```
$ echo '$$$' "* * * $HOME * * *" \ $HOME
$$$ * * * /homes/mgk25 * * * $HOME
```

The bash extension `$'...'` provides access to the full C string quoting syntax. For example `$'\x1b'` is the ASCII ESC character.

26

Job control (cont'd)

Foreground job: Stdin connected to terminal, shell prompt delayed until process exits, keyboard signals delivered to this single job.

Background job: Stdin disconnected (read attempt will suspend job), next shell prompt appears immediately, keyboard signals not delivered, shell prints notification when job terminates.

Keyboard signals: (keys can be changed with `stty` tool)

- ▶ Ctrl-C "intr" (SIGINT=2) by default aborts process
- ▶ Ctrl-\ "quit" (SIGQUIT=3) aborts process with core dump
- ▶ Ctrl-Z "susp" (SIGSTOP=19) suspends process

Another important signal (not available via keyboard):

- ▶ SIGKILL=9 destroys process immediately

28

Job control (cont'd)

Job control commands:

- ▶ `fg` resumes suspended job in foreground
- ▶ `bg` resumes suspended job in background
- ▶ `kill` sends signal to job or process

Job control commands accept as arguments

- ▶ process ID
- ▶ `%` + job number
- ▶ `%` + command name

Examples:

```
$ ghostview                # press Ctrl-Z
[6]+ Stopped                ghostview
$ bg
$ kill %6
```

29

Shell variables

Serve both as variables (of type string) in shell programming as well as environment variables for communication with programs.

Set *variable* to *value*:

```
variable=value
```

Note: No whitespace before or after “=” allowed.

Make variable visible to called programs:

```
export variable
export variable=value
```

Modify environment variables for one command only:

```
variable1=value variable2=value command
```

“`set`” shows all shell variables

“`printenv`” shows all (exported) environment variables.

31

Job control (cont'd)

A few more job control hints

- ▶ `kill -9 ...` sends SIGKILL to process. Should only be used as a last resort, if a normal kill (which sends SIGINT) failed, otherwise program has no chance to clean up resources before it terminates.
- ▶ The `jobs` command shows only jobs of the current shell, while `ps` and `top` list entire process table. Options for `ps` differ significantly between System V and BSD derivatives, check man pages.
- ▶ `fg %-` or just `%-` runs previously stopped job in foreground, which allows you to switch between several programs conveniently.

30

Standard environment variables

- ▶ `$HOME` — Your home directory, also available as “`~`”.
- ▶ `$USER/$LOGNAME` — Your login name.
- ▶ `$PATH` — Colon-separated list of directories in which shell looks for commands (e.g., “`/bin:/usr/bin:/usr/X11R6/bin`”).
Should never contain “`.`”, at least not at beginning. Why?
- ▶ `$LD_LIBRARY_PATH` — Colon-separated list of directories where the loader looks for shared libraries (see man `ld.so`)
- ▶ `$LANG, $LC_*` — Your “locale”, the name of a system-wide configuration file with information about your character set and language/country conventions (e.g., “`en_GB.UTF-8`”). `$LC_*` sets locale only for one category, e.g. `$LC_CTYPE` for character set and `$LC_COLLATE` for sorting order; `$LANG` sets default for everything. “`locale -a`” lists all available locales.
- ▶ `$TZ` — Specification of your timezone (mainly for remote users)
- ▶ `$OLDPWD` — Previous working directory, also available as “`~-`”.

32

Standard environment variables (cont'd)

- ▶ `$PS1` — The normal command prompt, e.g.

```
$ PS1='\[\033[7m\]\u@\h:\W \!\$ \[\033[m\] '
mgk25@ramsey:unixtools 12$
```
- ▶ `$PRINTER` — The default printer for `lpr`, `lpq` and `lprm`.
- ▶ `$TERM` — The terminal type (usually `xterm` or `vt100`).
- ▶ `$PAGER/$EDITOR` — The default pager/editor (usually `less` and `emacs`, respectively).
- ▶ `$DISPLAY` — The X server that X clients shall use.

man 7 environ

33

Plain-text files

- ▶ File is a sequence of lines (trad. each < 80 characters long).
- ▶ Characters ASCII encoded (or extension: ISO 8859-1 “Latin 1”, Microsoft’s CP1252, EUC, Unicode UTF-8, etc.)
- ▶ Each line ends with a special “line feed” control character (LF, Ctrl-J, byte value: $10_{10} = 0A_{16}$).
- ▶ “Horizontal tabulator” (HT, TAB, Ctrl-I, byte value: 9) advances the position of the next character to the next multiple-of-8 column.

Some systems (e.g., DOS, Windows, some Internet protocols) end each line instead with a **two** control-character sequence: “carriage return” (CR, Ctrl-M, $13_{10} = 0D_{16}$) plus “line feed”.

Different end-of-line conventions and different ASCII extensions make conversion of plain-text files necessary (`dos2unix`, `iconv`). Very annoying!

Alternative “flowed” plain-text format: no LF is stored inside a paragraph, line wrapping of paragraph happens only while a file is being displayed, LF terminates paragraph.

Some plain-text editors (e.g., Windows Notepad) start each UTF-8 plain-text file with a Unicode “Byte Order Mark” (BOM, $U+FEFF$, $EF_{16} BB_{16} BF_{16}$), which is not normally used under Unix.

35

Executable files and scripts

Many files signal their format in the first few “magic” bytes of the file content (e.g., `0x7f`, ‘E’, ‘L’, ‘F’ signals the System V *Executable and Linkable Format*, which is also used by Linux and Solaris).

The “file” tool identifies hundreds of file formats and some parameters based on a database of these “magic” bytes:

```
$ file $(which ls)
/bin/ls: ELF 32-bit LSB executable, Intel 80386
```

The kernel recognizes files starting with the magic bytes “#!” as “scripts” that are intended for processing by the interpreter named in the rest of the line, e.g. a bash script starts with

```
#!/bin/bash
```

If the kernel does not recognize a command file format, the shell will interpret each line of it, therefore, the “#!” is optional for shell scripts.

Use “`chmod +x file`” and “`./file`”, or “bash file”.

34

Shell compound commands

A *list* is a sequence of one or more pipelines separated by “;”, “&”, “&&” or “|”, and optionally terminated by one of “;”, “&” or end-of-line.

The return value of a list is that of the last command executed.

- ▶ `(list)` executes list in a subshell
- ▶ `{ list ; }` groups a list (to override operator priorities)
- ▶ `for variable in words ; do list ; done`

Expands *words* like command-line arguments, assigns one at a time to the *variable*, and executes *list* for each. Example:

```
for f in *.txt ; do cp $f $f.bak ; done
```

- ▶ `if list ; then list ; elif list ; then list ; else list ; fi`
- ▶ `while list ; do list ; done`
`until list ; do list ; done`

36

Shell compound commands (cont'd)

```
▶ case word in
    pattern|pattern|... ) list ;;
    ...
esac
```

Matches expanded *word* against each *pattern* in turn (same matching rules as pathname expansion) and executes the corresponding *list* when first match is found. Example:

```
case "$command" in
    start)
        app_server &
        processid=$! ;;
    stop)
        kill $processid ;;
    *)
        echo 'unknown command' ;;
esac
```

37

Boolean tests – examples

```
if [ -e $HOME/.rhosts ] ; then
    echo 'Found ~/.rhosts!' | \
    mail $LOGNAME -s 'Hacker backdoor?'
fi
```

Note: A backslash at the end of a command line causes end-of-line to be ignored.

```
if [ "`hostname`" == python.cl.cam.ac.uk ] ; then
    ( sleep 10 ; play ~/sounds/greeting.wav ) &
else
    xmessage 'Good Morning, Dave!' &
fi
[ "`arch`" != ix86 ] || { clear ; echo "I'm a PC" ; }
```

39

Boolean tests

The first *list* in the `if`, `while` and `until` commands is interpreted as a Boolean condition. The `true` and `false` commands return 0 and 1 respectively (note the inverse logic compared to Boolean values in C!).

The builtin command “`test expr`”, which can also be written as “[*expr*]” evaluates simple Boolean expressions on files, such as

```
-e file    is true if file exists.
-d file    is true if file exists and is a directory.
-f file    is true if file exists and is a normal file.
-r file    is true if file exists and is readable.
-w file    is true if file exists and is writable.
-x file    is true if file exists and is executable.
```

or strings, such as

```
string1 == string2    string1 < string2
string1 != string2    string1 > string2
```

38

Aliases and functions

Aliases allow a string to be substituted for the first word of a command:

```
$ alias dir='ls -la'
$ dir
```

Shell functions are defined with “`name () { list ; }`”. In the function body, the command-line arguments are available as `$1`, `$2`, `$3`, etc. The variable `$*` contains all arguments and `$#` their number.

```
$ unalias dir
$ dir () { ls -la $* ; }
```

Outside the body of a function definition, the variables `$*`, `$#`, `$1`, `$2`, `$3`, ... can be used to access the command-line arguments passed to a shell script.

40

Shell history

The shell records commands entered. These can be accessed in various ways to save keystrokes:

- ▶ “history” outputs all recently entered commands.
- ▶ “!n” is substituted by the n-th history entry.
- ▶ “!!” and “!-1” are equivalent to the previous command.
- ▶ “!*” is the previous command line minus the first word.
- ▶ Use cursor up/down keys to access history list, modify a previous command and reissue it by pressing Return.
- ▶ Type Ctrl-O instead of Return to issue command from history and edit its successor, which allows convenient repetition of entire command sequences.
- ▶ Type Ctrl-R to search string in history.

Most others probably only useful for teletype writers without cursor.

41

Readline

Interactive bash reads commands via the `readline` line-editor library. Many Emacs-like control key sequences are supported, such as:

- ▶ Ctrl-A/Ctrl-E moves cursor to **start/end** of line
- ▶ Ctrl-K deletes (**kills**) the rest of the line
- ▶ Ctrl-D **d**eletes the character under the cursor
- ▶ Ctrl-W deletes a **w**ord (first letter to cursor)
- ▶ Ctrl-Y inserts deleted strings
- ▶ ESC ^ performs history expansion on current line
- ▶ ESC # turns current line into a comment

Automatic word completion: Type the “Tab” key, and bash will complete the word you started when it is an existing \$variable, ~user, hostname, command or filename, depending on the context. If there is an ambiguity, pressing “Tab” a second time will show list of choices.

42

Startup files for terminal access

When you log in via a terminal line or telnet/rlogin/ssh:

- ▶ After verifying your password, the `login` command checks `/etc/passwd` to find out what shell to start for you.
- ▶ As a login shell, bash will execute the scripts
`/etc/profile`
`~/.profile`

The second one is where you can define your own environment. Use it to set exported variable values and trigger any activity that you want to happen at each login.

- ▶ Any subsequently started bash will read `~/.bashrc` instead, which is where you can define functions and aliases, which – unlike environment variables – are not exported to subshells.

43

Startup files for X Window System access

The “X server” provides access to display, keyboard and mouse for “X client” applications via the “X11 protocol”.

Before login, the only client is the X Display Manager (`xdm`).

After login, `xdm` will start the script `/usr/lib/X11/xdm/Xsession`. That invokes the “X clients” (`xterm`, etc.) that run on your desktop by default. If `~/.xsession` exists, this script will be called instead.

Most X clients in `Xsession` or `~/.xsession` are started in background, except for the last one, which is usually a window manager (`twm`, `fvwm2`, `KDE`, etc.). When this last client terminates, and with it the `Xsession` script, then `xdm` will reset the X server. This will terminate all X clients and the user is logged out.

You can configure your login screen in `~/.xsession`. You can also configure default parameters for many X clients via the `xrdb` command. See “man X” for details.

44

Typical .xsession file

```
#!/bin/bash
. ~/.profile

# set X defaults and keymaps
userresources=~/.Xdefaults
usermodmap=~/.Xmodmap
if [ -f $userresources ]; then
    xrdp $userresources
fi
if [ -f $usermodmap ]; then
    xmodmap $usermodmap
fi

# start some X clients as background processes
xterm -geometry 80x10+10+5 -C -title "`hostname -s` console" \
    -bg lightgreen &
xclock -geometry 80x80+0-0 -update 1 &
xload -geometry 80x80+90-0 -nolabel &

# start window manager as foreground process
if [ -x /usr/bin/X11/fvwm2 ]; then
    /usr/bin/X11/fvwm2
else
    twm
fi
```

45

sed – a stream editor

Designed to modify files in one pass and particularly suited for doing automated on-the-fly edits of text in pipes. sed scripts can be provided on the command line

```
sed [-e] 'command' files
```

or in a separate file

```
sed -f scriptfile files
```

General form of a sed command:

```
[address [,address]] [!] command [arguments]
```

Addresses can be line numbers or regular expressions. Last line is “\$”. One address selects a line, two addresses a line range (specifying start and end line). All commands are applied in sequence to each line. After this, the line is printed, unless option `-n` is used, in which case only the `p` command will print a line. The `!` negates address match. `{...}` can group commands per address.

46

sed – regular expressions

Regular expressions enclosed in `/.../`. Some regular expression meta characters:

- ▶ “.” matches any character (except new-line)
- ▶ “*” matches the preceding item zero or more times
- ▶ “+” matches the preceding item one or more times
- ▶ “?” matches the preceding item optionally (0–1 times)
- ▶ “^” matches start of line
- ▶ “\$” matches end of line
- ▶ “[...]” matches one of listed characters (use in character list “^” to negate and “-” for ranges)
- ▶ “\(...\)” grouping, “\{n,m\}” match n, \dots, m times
- ▶ “\” escape following meta character

47

sed – some examples

Substitute all occurrences of “Windows” with “Linux” (command: `s = substitute`, option: `g = “global” = all occurrences in line`):

```
sed 's/Windows/Linux/g'
```

Delete all lines that do not end with “OK” (command: `d = delete`):

```
sed '/OK$/!d'
```

Print only lines between those starting with BEGIN and END, inclusive:

```
sed -n '/^BEGIN/,/^END/p'
```

Substitute in lines 40–60 the first word starting with a capital letter with “X”:

```
sed '40,60s/[A-Z][a-zA-Z]*/X/'
```

48

grep, head, tail, sort

- ▶ Print only lines that contain pattern:

```
grep pattern files
```

Option `-v` negates match and `-i` makes match case insensitive.

- ▶ Print the first and the last 25 lines of a file:

```
head -n 25 file
```

```
tail -n 25 file
```

`tail -f` outputs growing file.

- ▶ Print the lines of a text file in alphabetical order: `sort file`
Options: `-k` select column, `-n` sort numbers, `-u` eliminate duplicate lines, `-r` reverse order.

49

find – traverse directory trees

`find directories expression` — recursively traverse the file trees rooted at the listed directories. Evaluate the Boolean expression for each file found. Examples:

Print relative pathname of each file below current directory:

```
$ find . -print
```

Erase each file named “core” below home directory if it was not modified in the last 10 days:

```
$ find ~ -name core -mtime +10 -exec rm -i {} \;
```

The test “`-mtime +10`” is true for files older than 10 days, concatenation of tests means “logical and”, so “`-exec`” will only be executed if all earlier terms were true. The “`{}`” is substituted with the current filename, and “`;`” terminates the list of arguments of the shell command provided to “`-exec`”.

51

chmod – set file permissions

- ▶ Unix file permissions: $3 \times 3 + 2 + 1 = 12$ bit information.
- ▶ Read/write/execute right for user/group/other.
- ▶ + set-user-id and set-group-id (elevated execution rights)
- ▶ + “sticky bit” (only owner can delete file from directory)
- ▶ `chmod ugoa[+|=]rwxst files`

Examples: Make file unreadable for anyone but the user/owner.

```
$ ls -l message.txt
-rw-r--r-- 1 mgk25 private 1527 Oct 8 01:05 message.txt
$ chmod go-rwx message.txt
$ ls -l message.txt
-rw----- 1 mgk25 private 1527 Oct 8 01:05 message.txt
```

For directories, “execution” right means right to traverse. Directories can be made traversable without being readable, such that only those who know the filenames inside can access them.

50

Some networking tools

- ▶ `wget url` — Fetch a file over the Internet via HTTP or FTP.
Option “`-r`” fetches HTML files recursively, option “`-l`” limits recursion depth.
- ▶ `ssh [user@]hostname [command]` — Log in via compressed and encrypted link to remote machine. If “`command`” is provided, execute it in remote shell, otherwise go interactive.
Preserves stdout/stderr distinction. Can also forward X11 requests (option “`-X`”) or arbitrary TCP/IP ports (options “`-L`” and “`-R`”) over secure link.
- ▶ `ssh-keygen -t dsa` — Generate DSA public/private key pair for password-free ssh authentication in “`~/.ssh/id_dsa.pub`” and “`~/.ssh/id_dsa`”. Protect “`id_dsa`” like a password!

Remote machine will not ask for password with ssh, if your private key “`~/.ssh/id_dsa`” fits one of the public keys (“locks”) listed on the remote machine in “`~/.ssh/authorized_keys`”.

On MCS Linux, your Novell-server home directory with `~/.ssh/authorized_keys` is mounted only **after** login, and therefore no password-free login for first session.

52

rsync

`rsync [options] source destination` — An improved cp.

- ▶ The source and/or destination file/directory names can be prefixed with `[user@]hostname`: if they are on a remote host.
- ▶ Uses `ssh` as a secure transport channel (may require `-e ssh`).
- ▶ Options to copy recursively entire subtrees (`-r`), preserve symbolic links (`-l`), permission bits (`-p`), and timestamps (`-t`).
- ▶ Will not transfer files (or parts of files) that are already present at the destination. An efficient algorithm determines, which bytes actually need to be transmitted only ⇒ very useful to keep huge file trees synchronised over slow links.

Application example: Very careful backup

```
rsync -e ssh -v -rlpt --delete --backup \  
  --backup-dir OLD/`date -Im` \  
  me@myhost.org:. mycopy/
```

Removes files at the destination that are no longer at the source, but keeps a timestamped copy of each changed or removed file in `mycopy/OLD/yyyy-mm-dd.../`, so nothing gets ever lost.

53

gzip & friends – compressing byte streams

- ▶ `gzip file` — convert “*file*” into a compressed “*file.gz*” (using a Lempel-Ziv/Huffman algorithm).
- ▶ `gunzip file` — decompress “*.gz” files.
- ▶ `[un]compress file` — [de]compress “*.Z” files (older tool using less efficient and patented LZW algorithm).
- ▶ `b[un]zip2 file` — [de]compress “*.bz2” files (newer tool using Burrows-Wheeler blocktransform).
- ▶ `zcat [file]` — decompress *.Z/*.gz to stdout for use in pipes.
- ▶ Extract compressed tar archive

```
$ zcat archive.tar.gz | tar xvf -  
$ tar xvzf archive.tgz          # GNU tar only!
```

55

tar – converting file trees into byte streams (and back)

Originally: “tape archiver”

Create archive (recurses into subdirectories):

```
$ tar cvf archive.tar files
```

Show archive content:

```
$ tar tvf archive.tar
```

Extract archive:

```
$ tar xvf archive.tar [files]
```

54

diff, patch – managing file differences

- ▶ `diff oldfile newfile` — Show difference between two text files, as lines that have to be inserted/deleted to change “*oldfile*” into “*newfile*”. Option “`-u`” gives better readable “unified” format with context lines. Option “`-r`” compares entire directory trees.

```
$ diff -u example.bak example.txt      $ diff example.bak example.txt  
--- example.bak                        2c2,3  
+++ example.txt                          < this sentence no verb  
@@ -1,2 +1,3 @@                          ---  
  an unmodified line                    > this sentence has a verb  
-this sentence no verb                    > a newly added line  
+this sentence has a verb  
+a newly added line
```

- ▶ `patch <diff-file` — Apply the changes listed in the provided diff output file to the old files named in it. The diff file should contain relative pathnames. If not, use option “`-pn`” to strip the first *n* directory names from pathnames in “*diff-file*”.

If the old files found by patch do not match exactly the removed lines in a “`-u`” diff output, patch will search whether the context lines can be located nearby and will report which line offset was necessary to apply them.

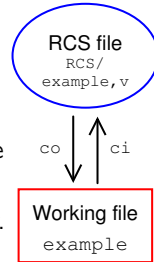
- ▶ `diff3 myfile oldfile yourfile` — Compare three files and merge the edits from different revision branches.

56

RCS – Revision Control System

Operates on individual files only. For every working file “example”, an associated RCS file “example,v” keeps a revision history database.

RCS files can be kept next to the working files, or in a subdirectory RCS/.



- ▶ `ci example` — Move a file (back) into the RCS file as the new latest revision (“check in”).
- ▶ `ci -u example` — Keep a read-only unlocked copy. “`ci -l...`” is equivalent to “`ci...`” followed by “`co...`”.
- ▶ `ci -l example` — Keep a writable locked copy (only one user can have the lock for a file at a time). “`ci -l...`” is equivalent to “`ci...`” followed by “`co -l...`”.
- ▶ `co example` — Fetches the latest revision from “example,v” as a read-only file (“check out”). Use option “`-rn.m`” to retrieve earlier revisions. There must not be a writable working file already.
- ▶ `co -l example` — Fetches the latest revision as a locked writable file if the lock is available.

57

RCS – Revision Control System (cont'd)

- ▶ `rcsdiff example` — Show differences between working file and latest version in repository (use option “`-rn.m`” to compare older revisions). Normal diff options like `-u` can be applied.
- ▶ `rlog example` — Show who made changes when on this file and left what change comments.

If you want to use RCS in a team, keep all the “*,v” files in a shared repository directory writable for everyone. Team members have their own respective working directory with a symbolic link named RCS to the shared directory.

As long as nobody touches the “*,v” files or manually changes the write permissions on working files, only one team member at a time can edit a file and old versions are never lost. The `rccs` command can be used by a team leader to bypass this policy and break locks or delete old revisions.

RCS remains useful for quickly maintaining history of single files, as an alternative to manually making backup copies of such files.

RCS is no longer commonly used for joint software development.

If you work in a distributed team on a project with subdirectories, need remote access, want to rename files easily, or simply hate locks, use `svn` or `git` instead.

58

svn – Subversion

Subversion is a popular centralized version control system (2001).

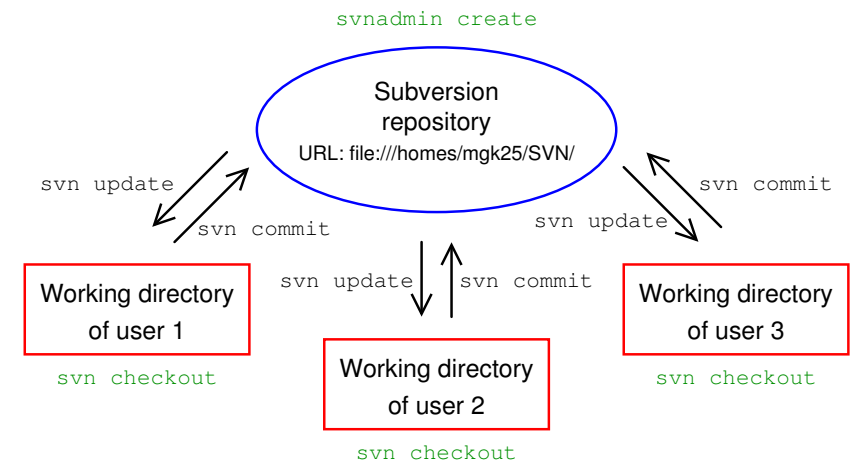
Main advantages over RCS:

- ▶ Supports a copy-modify-merge workflow (RCS: lock-modify-unlock). This allows team members to edit the same files **concurrently**.
 - Concurrent edits in different lines
⇒ merged automatically
 - Concurrent edits in the same lines
⇒ requires manual resolution of conflicts
- ▶ Manages entire directory trees, not just single files
- ▶ Understand tree edits (file copy, delete, rename, move)
- ▶ Supports several remote-access protocols (WebDAV, ssh, etc.)

Full documentation: <http://svnbook.red-bean.com/> and <http://subversion.apache.org/>
Microsoft Windows Subversion GUI: <http://tortoisesvn.tigris.org/>

59

svn – repository and working directories



Team administrator first creates repository: `svnadmin create`

Team members create personal working directories: `svn checkout`

Team members repeatedly fetch latest version: `svn update`

and return their changes: `svn commit`

60

svn – Subversion vs CVS

Subversion was specifically written to replace an older system, CVS, which in turn started out as a layer on top of RCS for managing entire directory trees. Its command-line interface closely follows that of CVS, but improves and simplifies the latter in many ways. In particular, Subversion

- ▶ understands renaming, moving, copying and replacing of both files and entire directory trees, no per-file version numbers
- ▶ understands symbolic links
- ▶ performs atomic commits
- ▶ versioned metadata (MIME types, EOL semantics, etc.)
- ▶ is easy to learn and understand for current CVS users
- ▶ simpler branching and tagging (through efficient copying)
- ▶ more efficient transactions, more disconnected operations
- ▶ wider choice of remote-access protocols (WebDAV, ssh, etc.)

Old CVS repositories can easily be converted: <http://cvs2svn.tigris.org/>

61

svn – directory edits

- ▶ `svn add filenames` — Put new files/folders under version control
Warning: adding a directory adds all content as well, unless you use `svn add -N dirnames`.
- ▶ `svn delete filenames` — Delete files/folders
- ▶ `svn copy source destination` — Copy files/folders
- ▶ `svn move source destination` — Move files/folders

The above four operations will not transfer the requested changes to the repository before the next `commit`, however the `delete/copy/move` operations perform the requested action immediately on your working files.

Remember not to use `rm/cp/mv` on working files that are under Subversion control, otherwise these operations will not be reflected in the repository after your next `commit`.

If you delete a version-controlled file with `rm`, the next `svn update` will restore it.

63

svn – setting up

Create new repository (e.g. `~/SVN/`), separate from working directory:

```
svnadmin create ~/SVN
```

Then checkout a working copy of the repository into a new working directory (`~/wdir`), referring to the repository via its URL:

```
svn checkout file://$HOME/SVN ~/wdir
```

Note that every subdirectory in your new working directory has a `.svn` subdirectory. This contains, among other things, the URL of your repository (see `svn info`). Therefore, inside the working directory, it is no longer necessary to add that repository URL as an argument to `svn` operations.

Now populate the repository with content:

- ▶ Create or move into the working directory some files
- ▶ Register them with Subversion (`svn add`)
- ▶ Push them into the repository (`svn commit`)

Then every team member, after their own initial `svn checkout`, does:

- ▶ Pull the latest version (`svn update`)
- ▶ Make edits
- ▶ Push them back to the repository (`svn commit`)

62

svn – querying the working-directory status

- ▶ `svn status` – List all files that differ between your working directory and the repository. The status code shown indicates:
 - A=added: this file will appear in the repository
 - D=deleted: this file will disappear from the repository
 - M=modified: you have edited this file
 - R=replaced: you used `svn delete` followed by `svn add`
 - C=conflict: at the last update, there was a conflict between your local changes and independent changes in the repository, which you still need to resolve manually
 - ?=unversioned: file is not in repository (suppress: `-q`)
 - !=missing: file in repository, but not in working dir.
- ▶ `svn diff [filenames]` — Show what you changed so far compared to the “base” version that you got at your last checkout or update.
- ▶ `svn info` — Show metadata about the working directory (revision of last update, URL of repository, etc.)

64

svn – commits and updates

- ▶ `svn commit [filenames]` — Check into the repository any modifications, additions, removals of files that you did since your last checkout or commit.

Option `-m '...'` provides a commit log message; without it, `svn commit` will call `$EDITOR` for you to enter one.

- ▶ `svn update [filenames]` — Apply modifications that others committed since you last updated your working directory.

This will list in the first column a letter for any file that differed between your working directory and the repository. Apart from the letter codes used by `status`, it also may indicate

- U=updated: get newer version of this file from repository
- G=merged: conflict, but was automatically resolved

Remaining conflicts (indicated as C) must be merged manually.

To assist in manual merging of conflicts, the update operation will write out all three file versions involved, all identified with appropriate filename extensions, as well as a diff3-style file that shows the differing lines next to each other for convenient editing.

svn – some more commands

- ▶ `svn resolved filenames` — Tell Subversion you have resolved a conflict. (Also cleans up the three additional files.)
- ▶ `svn revert filenames` — Undo local edits and go back to the version you had at your last checkout, commit, or update.
- ▶ `svn ls [filenames]` — List repository directory entries
- ▶ `svn cat filenames` — Output file contents from repository
Use `"svn cat filenames@rev"` to retrieve older revision `rev`.

Some of these commands can also be applied directly to a repository, without needing a working directory. In this case, specify a repository URL instead of a filename:

```
svn copy file://${HOME}/SVN/trunk \  
file://${HOME}/SVN/tags/release-1.0
```

An `svn copy` increases the repository size by only a trivial amount, independent of how much data was copied. Therefore, to give a particular version a symbolic name, simply `svn copy` it in the repository into a new subdirectory of that name.

65

66

Working example – User 1:

```
$ svnadmin create $HOME/svn-repos  
$ svn mkdir file://$HOME/svn-repos/example -m 'demo'  
Committed revision 1.  
$ svn checkout file://$HOME/svn-repos/example ex1  
Checked out revision 1.  
$ cd ex1  
ex1$ echo 'hello world' >file1  
ex1$ svn add file1  
A file1  
ex1$ svn commit -m 'adding my first file'  
Adding file1  
Transmitting file data .  
Committed revision 2.  
  
ex1$ echo 'bla' >file2  
A file2  
ex1$ svn add file2  
ex1$ echo 'hello humans' >file1  
ex1$ svn status  
M file1  
A file2  
ex1$ svn commit -m 'world -> humans'  
Sending file1  
Adding file2  
Transmitting file data ..  
Committed revision 3.  
ex1$ svn status  
  
ex1$ svn update  
U file1  
Updated to revision 4.  
ex2$ cat file1  
hello humans and dogs
```

Working example – User 2:

```
$ svn checkout file://$HOME/svn-repos/example ex2  
A ex2/file1  
Checked out revision 2.  
$ cd ex2  
  
ex2$ echo 'hello dogs' >file1  
ex2$ svn status  
M file1  
ex2$ svn commit -m 'world -> dogs'  
Sending file1  
svn: Commit failed (details follow):  
svn: File '/example/file1' is out of date  
ex2$ svn update  
Conflict discovered in 'file1'.  
Select: (p) postpone, (df) diff-full, (e) edit,  
(mc) mine-conflict, (tc) theirs-conflict,  
(s) show all options: p  
  
C file1  
A file2  
Updated to revision 3.  
Summary of conflicts:  
Text conflicts: 1  
ex2$ cat file1  
<<<<<< .mine  
hello dogs  
=====  
hello humans  
>>>>>> .r3  
ex2$ svn status  
? file1.mine  
? file1.r2  
? file1.r3  
C file1  
ex2$ echo 'hello humans and dogs' >file1  
ex2$ svn resolved file1  
Resolved conflicted state of 'file1'  
ex2$ svn status  
M file1  
ex2$ svn commit -m 'k9 extension'  
Sending file1  
Transmitting file data .  
Committed revision 4.
```

67

svn – remote access

The URL to an svn repository can point to a

- ▶ local file — `file://`
- ▶ Subversion/WebDAV Apache server — `http://` or `https://`
- ▶ Subversion server — `svn://`
- ▶ Subversion server accessed via ssh tunnel — `svn+ssh://`

The command

```
svn list svn+ssh://mgk25@linux2/home/mgk25/SVN/proj1
```

will ssh, as user `mgk25`, into host `linux2` and will start a server there with `svnserve -t`.

If you give others full shell access to your account to start `svnserve -t`, they could abuse this. Fortunately, `ssh` allows you to give others access to only a single program running under your user identity. You can add their public key to your `~/.ssh/authorized_keys` file with the option `command="..."` and other suitable restrictions (see `man sshd` and the `svn` book for details):

```
command="svnserve -t --tunnel-user=john -r /home/mgk25/SVN",no-port-forwarding,  
no-agent-forwarding,no-X11-forwarding,no-pty ssh-dss AAAB3...ogUc= john@bla.com
```

68

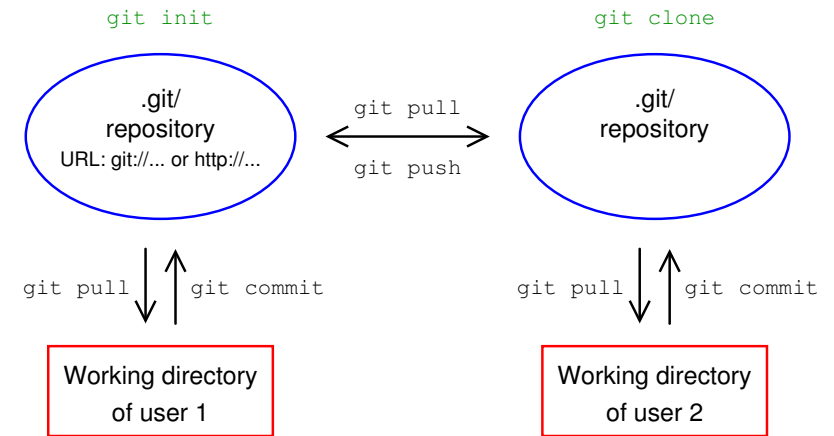
Distributed revision control systems

Popular examples: git, mercurial, bazaar

- ▶ No single central repository – more reliable and “democratic”
- ▶ Each participant holds a local repository of the revision history
- ▶ Committing and uploading a commit to another site are separate actions, i.e. commits, branches, etc. can be done off-line
- ▶ Creating and merging branches are quick and easy operations
- ▶ Branches do not have to be made visible to other developers
- ▶ Revisions identified by secure hash value rather than integer

Distributed version control systems are far more flexible and powerful than centralized ones (like Subversion), but require more experience.

git – Repository and working directories



Each working directory contains an associated repository in a `.git/` subdirectory of its top-level directory.

69

70

git – most basic commands

- ▶ `git init` – Create a new repository (inside a working directory). Option `--bare`, creates a stand-alone repository (for others to “push” into)
- ▶ `git clone` – Copy an existing repository and create an associated working directory around it (unless `--bare` is used). A cloned repository will remember the URL of the origin repository where it came from.
- ▶ `git pull` – Fetch updates from another repository (default: origin) and merge them into local working directory.
- ▶ `git add` – Add a new/modified file to the next commit.
- ▶ `git commit` – Save a revision to the local repository.
- ▶ `git push` – Forward local revisions to another repository.
- ▶ `git branch`, `git merge`, `git rebase`
- ▶ `gitk` – GUI history browser

Tutorial: <https://www.kernel.org/pub/software/scm/git/docs/gittutorial.html>

Manual: <https://www.kernel.org/pub/software/scm/git/docs/user-manual.html>

Git for computer scientists: <http://eagain.net/articles/git-for-computer-scientists/>

Git concepts simplified: <http://gitolite.com/gcs/>

Tutorial for Subversion users: <http://git.or.cz/course/svn.html>

71

Version-control etiquette

- ▶ **Use diff before commit:** carefully review what you actually changed. This often uncovers editing accidents and left-over temporary changes never intended to go into the repository.
- ▶ **Provide a useful commit message:** a complete, meaningful, honest, and accurate summary of everything you changed. Don't write just “bug fixed” (which?) or “API changed” (how and why?). Under git, there is a format convention for commit messages: a one-line summary, followed by an empty line, followed by details. Include function names, bug numbers.
- ▶ **Commit unrelated changes separately.** Others may later want to undo or merge them separately.
- ▶ **Commit related changes together.** Do not forget associated changes to documentation, test cases, build scripts, etc.
- ▶ **Leave the repository in a usable and consistent state.** It should always compile without errors and pass tests.
- ▶ **Avoid dependent or binary files in the repository.** Diffs on binary files are usually incomprehensible. Compiled output should be easy to recreate. It just wastes repository space and others can never be sure what is in it.

72

cc/gcc – the C compiler

Example:

```
$ cat hello.c
#include <stdio.h>
int main() { printf("Hello, World!\n"); return 0; }
$ gcc -o hello hello.c
$ ./hello
Hello, World!
```

Compiler accepts source (“*.c”) and object (“*.o”) files. Produces either final executable or object file (option “-c”). Common options:

- ▶ -W -Wall — activate warning messages (better analysis for suspicious code)
- ▶ -O — activate code optimizer
- ▶ -g — include debugging information (symbols, line numbers).

73

gdb – some common commands

- ▶ bt — print the current stack (backtracing function calls)
- ▶ p *expression* — print variable and expression values
- ▶ up/down — move between stack frames to inspect variables at different function call levels
- ▶ b ... — set breakpoint at specified line or function
- ▶ r ... — run program with specified command-line arguments
- ▶ s — continue until next source code line (skip function calls)
- ▶ n — continue until next source code line (follow function calls)

Also consider starting gdb within emacs with “ESC x gdb”, which causes the program-counter position to be indicated in source-file windows.

75

gdb – the C debugger

Best use on binaries compiled with “-g”.

- ▶ gdb binary — run command inside debugger (“r”) after setting breakpoints.
- ▶ gdb binary core — post mortem analysis on memory image of terminated process.

Enter in shell “ulimit -c 100000” before test run to enable core dumps. Core dump can be triggered by:

- ▶ a user pressing Ctrl-\ (SIGQUIT)
- ▶ a fatal processor or memory exception (segmentation violation, division by zero, etc.)

74

make – a project build tool

The files generated in a project fall into two categories:

- ▶ **Source files:** Files that cannot be regenerated easily, such as
 - working files directly created and edited by humans
 - files provided by outsiders
 - results of experiments
- ▶ **Derived files:** Files that can be recreated easily by merely executing a few shell commands, such as
 - object and executable code output from a compiler
 - output of document formatting tools
 - output of file-format conversion tools
 - results of post-processing steps for experimental data
 - source code generated by other programs
 - files downloaded from Internet archives

76

make – writing a Makefile

Many derived files have other source or derived files as *prerequisites*. They were generated from these input files and have to be regenerated as soon as one of the prerequisites has changed, and `make` does this.

A `Makefile` describes

- ▶ which (“target”) file in a project is derived
- ▶ on which other files that target depends as a prerequisite
- ▶ which shell command sequence will regenerate it

A `Makefile` contains rules of the form

```
target1 target2 ... : prereq1 prereq2 ...
    command1
    command2
    ...
```

Command lines **must** start with a TAB character (ASCII 9).

77

make – variables

Variables can be used to abbreviate rules:

```
CC=gcc
CFLAGS=-g -O
demo: demo.c demo.h
    $(CC) $(CFLAGS) -o $@ $<
```

```
data.gz: demo
    ./${< | gzip -c > $@
```

- ▶ `$@` — file name of the target of the rule
- ▶ `$<` — name of the first prerequisite
- ▶ `$(+)` — names of all prerequisites

Environment variables automatically become `make` variables, for example `$(HOME)`. A “\$” in a shell command has to be entered as “\$\$”.

79

make – writing a Makefile (cont’d)

Examples:

```
demo: demo.c demo.h
    gcc -g -O -o demo demo.c
```

```
data.gz: demo
    ./demo | gzip -c > data.gz
```

Call `make` with a list of target files as command-line arguments. It will check for every requested target whether it is still up-to-date and will regenerate it if not:

- ▶ It first checks recursively whether all prerequisites of a target are up to date.
- ▶ It then checks whether the target file exists and is newer than all its prerequisites.
- ▶ If not, it executes the regeneration commands specified.

Without arguments, `make` checks the targets of the first rule.

78

make – implicit rules, phony targets

Implicit rules apply to all files that match a pattern:

```
%.eps: %.gif
    giftopnm $< | pnmtops -noturn > $@
%.eps: %.jpg
    djpeg $< | pnmtops -noturn > $@
```

`make` knows a number of implicit rules by default, for instance

```
%.o: %.c
    $(CC) -c $(CPPFLAGS) $(CFLAGS) $<
```

It is customary to add rules with “phony targets” for routine tasks that will never produce the target file and just execute the commands:

```
clean:
    rm -f *~ *.bak *.o $(TARGETS) core
```

Common “phony targets” are “clean”, “test”, “install”.

80

perl – the Swiss Army Unix Tool

- ▶ a portable interpreted language with comprehensive library
- ▶ combines some of the features of C, sed, awk and the shell
- ▶ the expression and compound-statement syntax follows closely C, as do many standard library functions
- ▶ powerful regular expression and binary data conversion facilities make it well suited for parsing and converting file formats, extracting data, and formatting human-readable output
- ▶ offers arbitrary size strings, arrays and hash tables
- ▶ garbage collecting memory management
- ▶ dense and compact syntax leads to many potential pitfalls and has given Perl the reputation of a write-only hacker language
- ▶ widely believed to be less suited for beginners, numerical computation and large-scale software engineering, but highly popular for small to medium sized scripts, and Web CGI

81

perl – data types

Perl has three variable types, each with its own name space. The first character of each variable reference indicates the type accessed:

`$...` a scalar
`@...` an array of scalars
`%...` an associative array of scalars (hash table)

`[...]` selects an array element, `{...}` queries a hash table entry.

Examples of variable references:

`$days` = the value of the scalar variable "days"
`$days[28]` = element 29 of the array @days
`$days{'Feb'}` = the 'Feb' value from the hash table %days
`$#days` = last index of array @days
`@days` = (`$days[0]`, ..., `$days[$#days]`)
`@days[3,4,5]` = `@days[3..5]`
`@days{'a','c'}` = (`$days{'a'}`, `$days{'c'}`)
`%days` = (key1, val1, key2, val2, ...)

82

perl – scalar values

- ▶ A "scalar" variable can hold a string, number, or reference.
- ▶ Scalar variables can also hold the special `undef` value (set with `undef` and tested with `defined(...)`)
- ▶ Strings can consist of bytes or characters (Unicode/UTF-8).
More on Unicode character strings: `man perluniintro`.
- ▶ Numeric (decimal) and string values are automatically converted into each other as needed by operators.
(`5 - '3' == 2`, `'a' == 0`)
- ▶ In a Boolean context, the values `''`, `0`, `'0'`, or `undef` are interpreted as "false", everything else as "true".
Boolean operators return 0 or 1.
- ▶ References are typed pointers with reference counting.

83

perl – scalar literals

- ▶ Numeric constants follow the C format:
123 (decimal), 0173 (octal), 0x7b (hex), 3.14e9 (float)
Underscores can be added for legibility: `4_294_967_295`
- ▶ String constants enclosed with `"..."` will substitute variable references and other meta characters. In `'...'` only `'\'` and `'\"'` are substituted.

```
$header = "From: $name[$i]\@$host\n" .  
          "Subject: $subject{$msgid}\n";  
print 'Metacharacters include: $@%\\"';
```
- ▶ Strings can contain line feeds (multiple source-code lines).
- ▶ Multiline strings can also be entered with "here docs":

```
$header = <<"EOT";  
From: $name[$i]\@$host  
Subject: $subject{$msgid}  
EOT
```

84

perl – arrays

- ▶ Arrays start at index 0
- ▶ Index of last element of @foo is \$#foo (= length minus 1)
- ▶ Array variables evaluate in a scalar context to array length, i.e.
`scalar(@foo) == $#foo + 1;`
- ▶ List values are constructed by joining scalar values with comma operator (parenthesis often needed due to precedence rules):
`@foo = (3.1, 'h', $last);`
- ▶ Lists in lists lose their list identity: `(1, (2,3))` equals `(1,2,3)`
- ▶ Use `[...]` to generate reference to list (e.g., for nested lists).
- ▶ Null list: `()`
- ▶ List assignments: `($a,undef,$b,@c)=(1,2,3,4,5);` equals
`$a=1; $b=3; @c=(4,5);`
- ▶ Command line arguments are available in @ARGV.

85

perl – hash tables

- ▶ Literal of a hash table is a list of key/value pairs:
`%age = ('adam', 19, 'bob', 22, 'charlie', 7);`
Using => instead of comma between key and value increases readability:
`%age = ('adam' => 19, 'bob' => 22, 'charlie' => 7);`
- ▶ Access to hash table %age:
`$age{'john'} = $age{'adam'} + 6;`
- ▶ Remove entry: `delete $age{'charlie'};`
- ▶ Get list of all keys: `@family = keys %age;`
- ▶ Use `{...}` to generate reference to hash table.
- ▶ Environment variables are available in %ENV.

For more information: `man perldata`

86

perl – syntax

- ▶ Comments start with # and go to end of line (as in shell)
- ▶ Compound statements:

```
if (expr) block
  elsif (expr) block ...
  else block
while (expr) block [continue block]
for (expr; expr; expr) block
foreach var (list) block
```

Each *block* must be surrounded by {...} (no unbraced single statements as in C). The optional continue block is executed just before *expr* is evaluated again.
- ▶ The compound statements `if`, `unless`, `while`, and `until` can be appended to a statement:

```
$n = 0 if ++$n > 9;
do { $x >>= 1; } until $x < 64;
```

A do block is executed at least once.

87

perl – syntax (cont'd)

- ▶ Loop control:
 - `last` immediately exits a loop.
 - `next` executes the continue block of a loop, then jumps back to the top to test the expression.
 - `redo` restarts a loop block (without executing the continue block or evaluating the expression).
- ▶ The loop statements `while`, `for`, or `foreach` can be preceded by a label for reference in `next`, `last`, or `redo` instructions:

```
LINE: while (<STDIN>) {
  next LINE if /^#/;   # discard comments
  ...
}
```
- ▶ No need to declare global variables.

For more information: `man perlsyn`

88

perl – subroutines

- ▶ Subroutine declaration:

```
sub name block
```

- ▶ Subroutine call:

```
name(list);  
name list;  
&name;
```

A & prefix clarifies that a name identifies a subroutine. This is usually redundant thanks to a prior sub declaration or parenthesis. The third case passes @_ on as parameters.

- ▶ Parameters are passed as a flat list of scalars in the array @_.
- ▶ Perl subroutines are call-by-reference, that is \$_[0], ... are aliases for the actual parameters. Assignments to @_ elements will raise errors unless the corresponding parameters are lvalues.
- ▶ Subroutines return the value of the last expression evaluated or the argument of a return statement. It will be evaluated in the scalar/list context in which the subroutine was called.
- ▶ Use my(\$a,\$b); to declare local variables \$a and \$b within a block.

For more information: `man perlsub`

Example

```
sub max {  
    my ($x, $y) = @_;  
    return $x if $x > $y;  
    $y;  
}  
  
$m = max(5, 7);  
print "max = $m\n";
```

perl – operators

- ▶ Normal C/Java operators:

```
++ -- + - * / % << >> ! & | ^ && ||  
?: , = += -= *= ...
```

- ▶ Exponentiation: **

- ▶ Numeric comparison: == != <=> < > <= >=

- ▶ String comparison: eq ne cmp lt gt le ge

- ▶ String concatenation: \$a . \$a . \$a eq \$a x 3

- ▶ Apply regular expression operation to variable:
\$line =~ s/sed/perl/g;

- ▶ `...` executes a shell command

- ▶ .. returns list with a number range in a list context and works as a flip-flop in a scalar context (for sed-style line ranges)

For more information: `man perllop`

89

90

perl – references

Scalar variables can carry references to other scalar, list, hash-table, or subroutine values.

- ▶ To create a reference to another variable, subroutine or value, prefix it with \. (Much like & in C.)
- ▶ To dereference such a reference, prefix it with \$, @, %, or &, according to the resulting type. Use {...} around the reference to clarify operator precedence (\$\$a is short for \${\$a}).
- ▶ Hash-table and list references used in a lookup can also be dereferenced with ->, therefore \$a->{'john'} is short for \${\$a}{'john'} and \$b->[5] is short for \${\$b}[5].
- ▶ References to anonymous arrays can be created with [...].
- ▶ References to anonymous hash tables can be created with {...}.

For more information: `man perlref`

perl – examples of standard functions

`split /pattern/, expr`

Splits string into array of strings, separated by pattern.

`join expr, list`

Joins the strings in *list* into a single string, separated by value of *expr*.

`reverse list`

Reverse the order of elements in a list.

Can also be used to invert hash tables.

`substr expr, offset[, len]`

Extract substring.

Example:

```
$line = 'mgk25:x:1597:1597:Markus Kuhn:/homes/mgk25:/usr/bin/bash';  
@user = split(/:/, $line);  
($logname, $pw, $uid, $gid, $name, $home, $shell) = @user;  
$line = join(':', reverse(@user));
```

91

92

perl – more standard functions

<code>chop, chomp</code> Remove trailing character/linefeed from string	<code>lc, uc, lcfirst, ucfirst</code> Change entire string or first character to lowercase/uppercase
<code>pack, unpack</code> build/parse binary records	<code>chr, ord</code> ASCII ↔ integer conversion
<code>sprintf</code> format strings and numbers	<code>hex, oct</code> string → number conversion
<code>shift, unshift, push, pop</code> add/remove first/last array element	<code>wantarray</code> check scalar/list context in subroutine call
<code>die, warn</code> abort program with error/warning	<code>require, use</code> Import library module
<code>map, grep</code> Iterate over or filter list elements	

Perl provides most standard C and POSIX functions and system calls for arithmetic and low-level access to files, network sockets, and other interprocess communication facilities.

All built-in functions are listed in `man perlfunc`. A comprehensive set of add-on library modules is listed in `man perlmodlib` and thousands more are on <http://www.cpan.org/>.

93

perl – regular expressions

- ▶ Perl's regular expression syntax is similar to sed's, but `(){}` are metacharacters (and need no backslashes).
- ▶ Substrings matched by regular expression inside `(...)` are assigned to variables `$1, $2, $3, ...` and can be used in the replacement string of a `s/.../.../` expression.
- ▶ The substring matched by the regex pattern is assigned to `$&`, the unmatched prefix and suffix go into `$`` and `$'`.
- ▶ Predefined character classes include whitespace (`\s`), digits (`\d`), alphanumeric or `_` character (`\w`). The respective complement classes are defined by the corresponding uppercase letters, e.g. `\S` for non-whitespace characters.

Example:

```
$line = 'mgk25:x:1597:1597:Markus Kuhn:/homes/mgk25:/usr/bin/bash';
if ($line =~ /^(w+):[:]*:\d+:\d+:([:]*):[:]*:[^:]*$/) {
    $logname = $1; $name = $2;
    print "$logname = '$name'\n";
} else { die("Syntax error in '$line'\n"); }
```

For more information: `man perlre`

94

perl – predefined variables

`$_` The “default variable” for many operations, e.g.

```
print;           = print $_;
tr/a-z/A-Z/;     = $_ =~ tr/a-z/A-Z/;
while (<FILE>) ... = while ($_ = <FILE>) ...
```

`$.` Line number of the line most recently read from any file

`$?` Child process return value from the most recently closed pipe or ``...`` operator

`#!` Error message for the most recent system call, equivalent to C's `strerror(errno)`. Example:

```
open(FILE, 'test.dat') ||
die("Can't read 'test.dat': $!\n");
```

For many more: `man perlvar`

95

perl – file input/output

- ▶ `open filehandle, expr`

```
open(F1, 'test.dat'); # open file 'test.dat' for reading
open(F2, '>test.dat'); # create file 'test.dat' for writing
open(F3, '>>test.dat'); # append to file 'test.dat'
open(F4, 'date|'); # invoke 'date' and connect to its stdout
open(F5, '|mail -s test'); # invoke 'mail' and connect to its stdin
```
- ▶ `print filehandle, list`
- ▶ `close, eof, getc, seek, read, format, write, truncate`
- ▶ “`<filehandle>`” reads another line from file handle FILE and returns the string. Used without assignment in a while loop, the line read will be assigned to `$_`.
- ▶ “`<>`” opens one file after another listed on the command line (or stdin if none given) and reads out one line each time.

96

perl – invocation

- ▶ First line of a Perl script: `#!/usr/bin/perl` (as with shell)
- ▶ Option “-e” reads code from command line (as with sed)
- ▶ Option “-w” prints warnings about dubious-looking code.
- ▶ Option “-d” activates the Perl debugger (see `man perldebug`)
- ▶ Option “-p” places the loop

```
while (<>) { ... print; }
```

around the script, such that `perl` reads and prints every line. This way, Perl can be used much like `sed`:

```
sed -e 's/sed/perl/g'  
perl -pe 's/sed/perl/g'
```
- ▶ Option `-n` is like `-p` without the “`print;`”.
- ▶ Option “-i[*backup-suffix*]” adds in-place file modification to `-p`. It renames the input file, opens an output file with the original name and directs the input into it.

For more information: `man perlrun`

97

perl – a simple example

Generate a list of email addresses of everyone on the Computer Lab's “People” web page, sorted by surname.

Example input:

```
...  
<tr><td><a NAME="asa28">asa28</a></td><td>FE04</td><td>63622</td><td></td></tr>  
></td><td><a HREF="/users/asa28/">Abrahams, Alan</a></td></tr>  
<tr><td><a NAME="mha23">mha23</a></td><td>FE22</td><td>63692</td><td></td></tr>  
></td><td><a HREF="/users/mha23/">Allen-Williams, Mair</a></td></tr>  
<tr><td><a NAME="sa333">sa333</a></td><td>GC33</td><td>63680</td><td></td></tr>  
></td><td>Allott, Stephen</td></tr>  
...
```

Example output:

```
Alan Abrahams <asa28@cl.cam.ac.uk>  
Mair Allen-Williams <mha23@cl.cam.ac.uk>  
Stephen Allott <sa333@cl.cam.ac.uk>
```

99

perl – a stream editing example

“Spammers” send out unsolicited bulk email (junk email), for marketing or fraud. They harvest millions of valid email addresses, for example from web pages.

To make email addresses in your web pages slightly harder to harvest, you can avoid the “@” sign in the HTML source. For instance, convert

```
<a href="mailto:jd@acm.org">jd@acm.org</a>
```

into

```
<a href="mailto:jd%40acm.org">jd%#64;acm.org</a>
```

The lines

```
perl -pi.bak <<EOT *.html  
s/(href="mailto:[^@\\"]+)(@[^@\\"]+\\")/$1%40$2/ig;  
s/([a-zA-Z0-9\\.\\-\\+\\_]+)(@[a-zA-Z0-9\\.\\-\\+\\_]+)/$1&#64;$2/ig;  
EOT
```

will do that transformation for all HTML files in your directory.

More sophisticated methods: hide an email address in JavaScript or use a CAPTCHA.

98

perl – a simple example

Possible solution:

```
#!/usr/bin/perl  
$url = 'http://www.cl.cam.ac.uk/UoCCL/people/directory.html';  
open(HTML, "wget -O - '$url' |") || die("Can't start 'wget': $!\n");  
while (<HTML>) {  
    if (/^<tr><td><a name="(\\w+)">.*</tr>$/i) {  
        $csrid = $1;  
        if (/<td><a href="["]*">?([<>]*), ([<>]*)(</a>)?</td></tr>$/i) {  
            $email{$csrid} = "$3 $2 <$csrid@cl.cam.ac.uk>";  
            $surname{$csrid} = $2;  
        } else { die ("Syntax error:\n$_") }  
    }  
}  
foreach $s (sort({$surname{$a} cmp $surname{$b}} keys(%email))) {  
    print "$email{$s}\n";  
}
```

Warning: This simple-minded solution makes numerous assumptions about how the web page is formatted, which may or may not be valid. Can you name examples of what could go wrong?

100

perl – email-header parsing example

Email headers (as defined in RFC 822) have the form:

```
$header = <<'EOT';
From Ian.Grant@c1.cam.ac.uk  21 Sep 2004 10:10:18 +0100
Received: from ppsw-8.csi.cam.ac.uk ([131.111.8.138])
        by mta1.c1.cam.ac.uk with esmtp (Exim 3.092 #1)
        id 1V9afA-0004E1-00 for Markus.Kuhn@c1.cam.ac.uk;
        Tue, 21 Sep 2004 10:10:16 +0100
Date: Tue, 21 Sep 2004 10:10:05 +0100
To: Markus.Kuhn@c1.cam.ac.uk
Subject: Re: Unix tools notes
Message-ID: <514FGFED.mail1VJ3982Y@c1.cam.ac.uk>
EOT
```

This can be converted into a Perl hash table as easily as

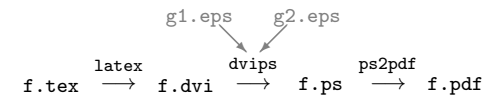
```
$header =~ s/\n\s+/ /g; # fix continuation lines
%hdr    = (FROM => split /\s*/m, $header);
```

and accessed as in `if ($hdr{Subject} =~ /Unix tools/) ...`

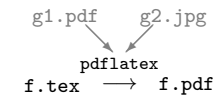
L^AT_EX – a document formatter

L^AT_EX is a sophisticated macro package for the T_EX text formatting system. Thanks to its excellent facilities for mathematical typesetting, it has become the de-facto standard for preparing scientific publications in mathematical, physical, computing and engineering disciplines.

Classic processing steps:



Modern alternative:

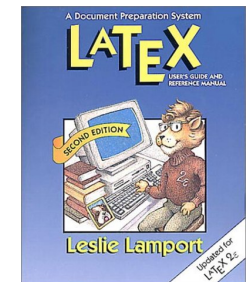


Recommended introduction:

Leslie Lamport: L^AT_EX – a document preparation system. 2nd ed., Addison-Wesley, 1994.

Online tutorials: <http://www.latex-project.org/guides/>
T_EX Frequently Asked Questions: <http://www.tex.ac.uk/cgi-bin/textfaq2html>

For advanced users: Mittelbach, et al.: *The L^AT_EX Companion*. 2nd ed., Addison-Wesley, 2004.



101

102

L^AT_EX example

```
\documentclass[12pt]{article}
\setlength{\textwidth}{75mm}
\begin{document}
\title{\TeX -- a summary}
\author{Markus Kuhn}
\date{28 November 2013}
\maketitle
\thispagestyle{empty}

\section{Introduction}

Mathematical formulæ such as
 $e^{i\pi} = -1$  or even
 $\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-\frac{1}{2}x^2}$ 
were a real 'pain' to typeset until
\textsc{Knuth}'s text formatter \TeX
became available \cite{Knuth86}.

\begin{thebibliography}{9}
\bibitem{Knuth86}Donald E. Knuth:
The \TeX book. Addison-Wesley, 1986.
\end{thebibliography}

\end{document}
```

T_EX – a summary

Markus Kuhn

28 November 2013

1 Introduction

Mathematical formulæ such as $e^{i\pi} = -1$ or even

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-\frac{1}{2}x^2}$$

were a real 'pain' to typeset until KNUTH's text formatter T_EX became available [1].

References

[1] Donald E. Knuth: *The T_EXbook*. Addison-Wesley, 1986.

103

T_EX input syntax

- ▶ T_EX reads plain-text *.tex files (e.g., prepared with emacs)
- ▶ no distinction is made between space character and line feed
- ▶ multiple spaces are treated like a single space
- ▶ multiple line feeds (empty lines) are treated as a paragraph separator (just like the \par command)
- ▶ command, macro and variable names start with a backslash (\), followed by either a sequence of letters or a single non-letter character (uppercase/lowercase is significant).
Correct: \par, \item, \pagethree, \LaTeX, \+, \\\, \3
Wrong: \page33, \<>
- ▶ space and line-feed characters are ignored if they follow a command/macro/variable name consisting of letters. Use _ to add an explicit space (e.g., \TeX\ syntax ⇒ T_EX syntax).

104

Characters with special semantics

In *.tex input files, the characters

\$ % & ~ _ ^ \ { }

have special functions. Some of these can be included in regular text by writing

\# \\$ \% \& _ \^ \{ \}

L^AT_EX supports typesetting all ASCII characters via the \verb and \url macros.

% starts a comment

All characters between (and including) a % and the next line feed will be ignored. Append % at the end of a line to avoid interpretation of the subsequent line feed as a space.

[# plus a digit denotes a parameter in macros, ~ is a no-break space, \$ delimits inline equations, & is used as a tabulator mark, \ is a line separator, ^ indicates a superscript and _ a subscript in math mode.]

105

Blocks

State changes inside a { ... } block last only until the next }:

{This is a \bf bold} statement.

↓

This is a **bold** statement.

Commands and macros read for each argument either a single character or a block enclosed by { and }:

Typeset \textsl M in \textsl{slanted style}.

↓

Typeset *M* in *slanted style*.

Values of optional L^AT_EX macro arguments are enclosed by [...].

106

Typewriting versus Typesetting

The ASCII (ISO 646) 7-bit character set with its 94 graphic characters

!"#\$%&'()*+,-./0123456789:;<=>?
@ABCDEFGHIJKLMNQRSTUvwxyz[\]^_
`abcdefghijklmnopqrstuvwxyz{|}~

was designed to cover the character repertoire of US typewriters and teletype printers. Some new symbols such as [\]{|}_ were added in the hope that they will be useful for programming.

T_EX defines a number of shortcuts and macros to access the full range of “typographic” characters used in high-quality book printing. These still cannot be found on the standard PC keyboard, which was designed for 7-bit ASCII.

107

Dashes

ASCII provides only a single combined hyphen-minus character, but typesetters distinguish carefully between several dash characters:

-	⇒	-	hyphen
--	⇒	-	en dash
---	⇒	—	em dash
\$\$-	⇒	-	minus

The hyphen (-) is the shortest of these and is used to combine separate words or split words across line-breaks.

The en dash (–) is often used to denote a range of numbers (as in pages 64–128), or – as in this example – as a punctuation dash.

The em dash is used—like this—as a punctuation dash, often without surrounding space, especially in US typography.

The minus (−) is a mathematical operator, whose shape matches the plus (+), unlike the hyphen or dashes. Compare: −+, −+, —+, −+.

108

Quotation marks

Typewriters and ASCII offer only unidirectional 'single' and "double" quotation marks, while typesetters use ‘curly’ and “directed” variants.

T_EX input files use the single quotation mark (') and the grave accent (̀) to encode these, as well the mathematical ‘prime’ marker and the French accents:

˘	⇒	‘	left quote
'	⇒	’	right quote
˘˘	⇒	“	left doublequote
''	⇒	”	right doublequote
\$'\$	⇒	'	prime
\'u	⇒	ú	acute accent
\`u	⇒	ù	grave accent

The apostrophe (it's) is identical to the right single quotation mark.

In some older terminal fonts (especially of US origin), the ˘ and ' characters have a compromise shape somewhere between the quotation marks ‘’ and the accents ˘˘.

Non-ASCII Symbols

ı	!˘	Å	\AA	¶	\P
ı	?˘	ø	\o	†	\dag
œ	\oe	Ø	\O	‡	\ddag
Œ	\OE	†	\l	©	\copyright
æ	\ae	Ł	\L	£	\pounds
Æ	\AE	ß	\ss	...	\ldots
å	\aa	§	\S		

Combining characters

ó	\'o	ō	\=o	ô	\t{oo}
ò	\`o	ó	\.o	o	\c{o}
ô	\^o	ö	\u{o}	o	\d{o}
ö	\"o	õ	\v{o}	o	\b{o}
õ	\~o	ö	\H{o}		

109

110

Space – the final frontier

Traditional English typesetting inserts a larger space at the end of a sentence. T_EX believes any space after a period terminates a sentence, unless it is preceded by an uppercase letter. Parenthesis are ignored.

This works often: J. F. Kennedy's U.S. budget. Look!

But not always: E.g. NASA. Dr. K. Smith et al. agree.

To correct failures of this heuristic, use

~	⇒	no-break space
_	⇒	force normal space
\@	⇒	following punctuation ends sentence

as in

E.g. \ NASA \@. Dr. ~K. Smith et al. \ agree.

↓
E.g. NASA. Dr. K. Smith et al. agree.

Or disable the distinction of spaces with \frenchspacing.

Structure of a L^AT_EX document

First select a document class and its options, e.g. with

```
\documentclass[12pt,a4paper]{article}
```

Standard classes: article, report, book, letter, slides.

Publishers often provide authors with their own class as a *.cls file. Appendix A of *The L^AT_EX Companion* explains how to write new class files. A popular class for presentation slides: beamer

Delimit block environments as in

```
\begin{document} ... \end{document}
```

Others: abstract, center, verbatim, itemize, tabular, ...

Mark headings with

```
\section{...}           \subsection{...}
\subsubsection{...}    \paragraph{...}
```

and L^AT_EX will take care of font sizes, numbering, and table of contents.

T_EX is a full programming language with macros, variables, recursion, conditional branching, file I/O, and a huge collection of add-ons.

111

112

Tweaking and extending L^AT_EX

L^AT_EX behaviour can be changed by overwriting predefined variables and macros. This can be done

- ▶ in the *preamble* (between the `\documentclass{...}` and `\begin{document}` lines) \Rightarrow for the entire document
- ▶ anywhere in the document \Rightarrow the effect will last only until the end of the current block (i.e., the next `}` or `\end{...}`)

Packages

A huge collection of extension packages exists for L^AT_EX. Some merely define additional macros and environments, others rewrite parts of L^AT_EX's internal machinery. For example, adding to the preamble

```
\usepackage{hyperref}
```

loads all the macros and settings defined in the `hyperref.sty` package.

`hyperref` adds new macros, such as `\url{...}` for typesetting URLs, but also automatically turns every reference to a page, section, or bibliographic entry into a hyperlink.

Documentation: `texdoc packagename` e.g. `texdoc geometry`

113

Example: changing page layout geometry

Adjust margins manually, via numerous length variables:

```
\setlength{\oddsidemargin}{-0.4mm} % 25 mm left margin
\setlength{\evensidemargin}{\oddsidemargin}
\setlength{\textwidth}{160mm} % 25 mm right margin
\setlength{\topmargin}{-5.4mm} % 20 mm top margin
\setlength{\headheight}{5mm}
\setlength{\headsep}{5mm}
\setlength{\footskip}{10mm}
\setlength{\textheight}{237mm} % 20 mm bottom margin
```

More comfortable:

```
\usepackage[vmargin=20mm,hmargin=25mm]{geometry}
```

The `geometry.sty` package automatically recalculates any dimensions not specified.

Make paragraphs not indented at the first line, but spaced apart slightly:

```
\setlength{\parindent}{0mm}
\setlength{\parskip}{\medskipamount}
```

Or just:

```
\usepackage{parskip}
```

114

Mathematical typesetting

In T_EX, mathematical formulas are formatted in a completely different mode from that used for normal text.

Inline formulas such as a_n (`a_n`) that appear as part of a normal paragraph have to be surrounded with `$...$`, while displayed formulas such as

$$F_n = F_{n-1} + F_{n-2} \quad (\backslash[F_n=F_{n-1}+F_{n-2}])$$

are entered in between `\[...]`. In math mode

- ▶ space characters are ignored; T_EX adds its own space around operators based on heuristics; manually add `thinspace` with `"\,`
- ▶ a special math italic font with different inter-character spacing is used, to show single-letter variables better in products
- ▶ many additional macros for special symbols are defined

Math italic is very *different* and not suitable for writing words! Use

```
\mathrm{...} around words, as in  $v_{\mathrm{diff}} \rightarrow v_{\text{diff}}$ .
```

Macros for neatly aligning multiple equations: `\usepackage{amsmath}`, see `texdoc amsldoc`.

115

Mathematical symbols

Greek letters

Γ	<code>\Gamma</code>	δ	<code>\delta</code>	π	<code>\pi</code>
Δ	<code>\Delta</code>	ϵ	<code>\epsilon</code>	ϖ	<code>\varpi</code>
Θ	<code>\Theta</code>	ε	<code>\varepsilon</code>	ρ	<code>\rho</code>
Λ	<code>\Lambda</code>	ζ	<code>\zeta</code>	ϱ	<code>\varrho</code>
Ξ	<code>\Xi</code>	η	<code>\eta</code>	σ	<code>\sigma</code>
Π	<code>\Pi</code>	θ	<code>\theta</code>	ς	<code>\varsigma</code>
Σ	<code>\Sigma</code>	ϑ	<code>\vartheta</code>	τ	<code>\tau</code>
Υ	<code>\Upsilon</code>	ι	<code>\iota</code>	υ	<code>\upsilon</code>
Φ	<code>\Phi</code>	κ	<code>\kappa</code>	ϕ	<code>\phi</code>
Ψ	<code>\Psi</code>	λ	<code>\lambda</code>	φ	<code>\varphi</code>
Ω	<code>\Omega</code>	μ	<code>\mu</code>	χ	<code>\chi</code>
α	<code>\alpha</code>	ν	<code>\nu</code>	ψ	<code>\psi</code>
β	<code>\beta</code>	ξ	<code>\xi</code>	ω	<code>\omega</code>
γ	<code>\gamma</code>	o	<code>o</code>		

116

Mathematical symbols

Binary operators

\pm	<code>\pm</code>	\triangleleft	<code>\lhd</code>	\oplus	<code>\oplus</code>
\mp	<code>\mp</code>	\cap	<code>\cap</code>	\ominus	<code>\ominus</code>
\setminus	<code>\setminus</code>	\cup	<code>\cup</code>	\otimes	<code>\otimes</code>
\cdot	<code>\cdot</code>	\uplus	<code>\uplus</code>	\oslash	<code>\oslash</code>
\times	<code>\times</code>	\sqcap	<code>\sqcap</code>	\odot	<code>\odot</code>
\ast	<code>\ast</code>	\sqcup	<code>\sqcup</code>	\dagger	<code>\dagger</code>
\star	<code>\star</code>	\wr	<code>\wr</code>	\ddagger	<code>\ddagger</code>
\diamond	<code>\diamond</code>	\bigcirc	<code>\bigcirc</code>	\amalg	<code>\amalg</code>
\circ	<code>\circ</code>	\triangleright	<code>\rhd</code>	\triangleleft	<code>\unlhd</code>
\bullet	<code>\bullet</code>	\vee	<code>\vee</code>	\triangleright	<code>\unrhd</code>
\div	<code>\div</code>	\wedge	<code>\wedge</code>		
\triangleleft	<code>\triangleleft</code>	\triangleup	<code>\bigtriangleup</code>		
\triangleright	<code>\triangleright</code>	∇	<code>\bigtriangledown</code>		

117

Mathematical symbols

Relations

\leq	<code>\leq</code>	\geq	<code>\geq</code>	\equiv	<code>\equiv</code>
\prec	<code>\prec</code>	\succ	<code>\succ</code>	\sim	<code>\sim</code>
\preceq	<code>\preceq</code>	\succeq	<code>\succeq</code>	\simeq	<code>\simeq</code>
\ll	<code>\ll</code>	\gg	<code>\gg</code>	\asymp	<code>\asymp</code>
\subset	<code>\subset</code>	\supset	<code>\supset</code>	\approx	<code>\approx</code>
\subseteq	<code>\subseteq</code>	\supseteq	<code>\supseteq</code>	\cong	<code>\cong</code>
\sqsubset	<code>\sqsubset</code>	\sqsupseteq	<code>\sqsupseteq</code>	\bowtie	<code>\bowtie</code>
\in	<code>\in</code>	\ni	<code>\ni</code>	\propto	<code>\propto</code>
\vdash	<code>\vdash</code>	\dashv	<code>\dashv</code>	\models	<code>\models</code>
\smile	<code>\smile</code>	\mid	<code>\mid</code>	\doteq	<code>\doteq</code>
\frown	<code>\frown</code>	\parallel	<code>\parallel</code>	\perp	<code>\perp</code>
\sqsubset	<code>\sqsubset</code>	\sqsupset	<code>\sqsupset</code>	\Join	<code>\Join</code>
$\not<$	<code>\not<</code>	$\not=$	<code>\not=</code>	$\not>$	<code>\not></code>
$\not\leq$	<code>\not\leq</code>	$\not\geq$	<code>\not\geq</code>	$\not\equiv$	<code>\not\equiv</code>
$\not\prec$	<code>\not\prec</code>	$\not\succ$	<code>\not\succ</code>	\dots	<code>\dots</code>

118

Mathematical symbols

Arrows

\leftarrow	<code>\leftarrow</code>	\longleftrightarrow	<code>\Longleftrightarrow</code>
\Lleftarrow	<code>\Lleftarrow</code>	\longmapsto	<code>\longmapsto</code>
\rightarrow	<code>\rightarrow</code>	\hookrightarrow	<code>\hookrightarrow</code>
\Rightarrow	<code>\Rightarrow</code>	\rightharpoonup	<code>\rightharpoonup</code>
\leftrightarrow	<code>\leftrightarrow</code>	\rightharpoondown	<code>\rightharpoondown</code>
\Lleftrightarrow	<code>\Lleftrightarrow</code>	\leadsto	<code>\leadsto</code>
\mapsto	<code>\mapsto</code>	\Uparrow	<code>\Uparrow</code>
\hookleftarrow	<code>\hookleftarrow</code>	\Uparrow	<code>\Uparrow</code>
\leftharpoonup	<code>\leftharpoonup</code>	\Downarrow	<code>\Downarrow</code>
\leftharpoondown	<code>\leftharpoondown</code>	\Updownarrow	<code>\Updownarrow</code>
\rightleftharpoons	<code>\rightleftharpoons</code>	\Updownarrow	<code>\Updownarrow</code>
\longleftarrow	<code>\longleftarrow</code>	\nearrow	<code>\nearrow</code>
\Llongleftarrow	<code>\Llongleftarrow</code>	\searrow	<code>\searrow</code>
\longrightarrow	<code>\longrightarrow</code>	\swarrow	<code>\swarrow</code>
\Longrightarrow	<code>\Longrightarrow</code>	\nwarrow	<code>\nwarrow</code>
\longleftrightarrow	<code>\longleftrightarrow</code>		

119

Mathematical symbols

\aleph	<code>\aleph</code>	\prime	<code>\prime</code>	\forall	<code>\forall</code>
\hbar	<code>\hbar</code>	\emptyset	<code>\emptyset</code>	\exists	<code>\exists</code>
\imath	<code>\imath</code>	∇	<code>\nabla</code>	\neg	<code>\neg</code>
\jmath	<code>\jmath</code>	\surd	<code>\surd</code>	\flat	<code>\flat</code>
ℓ	<code>\ell</code>	\top	<code>\top</code>	\natural	<code>\natural</code>
\wp	<code>\wp</code>	\bot	<code>\bot</code>	\sharp	<code>\sharp</code>
\Re	<code>\Re</code>	\parallel	<code>\parallel</code>	\clubsuit	<code>\clubsuit</code>
\Im	<code>\Im</code>	\angle	<code>\angle</code>	\diamondsuit	<code>\diamondsuit</code>
∂	<code>\partial</code>	\triangle	<code>\triangle</code>	\heartsuit	<code>\heartsuit</code>
∞	<code>\infty</code>	\backslash	<code>\backslash</code>	\spadesuit	<code>\spadesuit</code>
\Box	<code>\Box</code>	\diamond	<code>\Diamond</code>		
\dots	<code>\ldots</code>	\cdots	<code>\cdots</code>	\vdots	<code>\vdots</code>
		\ddots	<code>\ddots</code>		

120

Mathematical symbols

Large operators

\sum	<code>\sum</code>	\bigcap	<code>\bigcap</code>	\bigodot	<code>\bigodot</code>
\prod	<code>\prod</code>	\bigcup	<code>\bigcup</code>	\bigotimes	<code>\bigotimes</code>
\coprod	<code>\coprod</code>	\bigsqcup	<code>\bigsqcup</code>	\bigoplus	<code>\bigoplus</code>
\int	<code>\int</code>	\bigvee	<code>\bigvee</code>	\biguplus	<code>\biguplus</code>
\oint	<code>\oint</code>	\bigwedge	<code>\bigwedge</code>		

Delimiters

$[$	<code>\lbrack</code>	$]$	<code>\rbrack</code>
\lfloor	<code>\lfloor</code>	\rfloor	<code>\rfloor</code>
\lceil	<code>\lceil</code>	\rceil	<code>\rceil</code>
$\{$	<code>\lbrace</code>	$\}$	<code>\rbrace</code>
\langle	<code>\langle</code>	\rangle	<code>\rangle</code>
\llbracket	<code>\llbracket</code>	\rrbracket	<code>\rrbracket</code>
$\langle\!\langle$	<code>\langle\!\langle</code>	$\rangle\!\rangle$	<code>\rangle\!\rangle</code>

121

Including graphics

DVI only supports characters and filled rectangles, but dvips and pdftex also understand embedded “special” instructions that provide more.

Embedded PostScript (EPS) vector graphics:

Normal PostScript files (*.ps) produce a sequence of pages. An EPS file describes only an image and is meant to be included into a PostScript page. EPS files lack instructions to output paper, but define a rectangular “bounding box”, using special `%%BoundingBox:` comments.

Load the `graphicx` extension of L^AT_EX by adding

```
\usepackage{graphicx}
```

to the preamble. Then write

```
\includegraphics{filename.eps}
```

wherever you want to include the graphics file into your text.

In `pdflatex`, the `graphicx` package allows you to include graphics from PDF (vector graphics), JPEG (photos) and PNG (bitmap) files:

```
\includegraphics{filename.pdf}
```

123

Mathematical symbols

Alternative names

\neq	<code>\ne</code>	$\{$	<code>\{</code>	\ni	<code>\owns</code>	$ $	<code>\vert</code>
\neq	<code>\neq</code>	$\}$	<code>\}</code>	\wedge	<code>\land</code>	$\ $	<code>\Vert</code>
\leq	<code>\le</code>	\rightarrow	<code>\to</code>	\vee	<code>\lor</code>		
\geq	<code>\ge</code>	\leftarrow	<code>\gets</code>	\neg	<code>\lnot</code>		

Stacking things

a^b	<code>a^{b}</code>	a_b	<code>a_{b}</code>
$\overline{a-b}$	<code>\overline{a-b}</code>	$\overbrace{a-b}$	<code>\overbrace{a-b}</code>
$\underline{a-b}$	<code>\underline{a-b}</code>	$\underbrace{a-b}$	<code>\underbrace{a-b}</code>
$= \begin{cases} a^{2^2}, & a \geq 0 \\ -a, & a < 0 \end{cases}$	<code>=\left\{\begin{array}{c} a^{2^2}, & a \ge 0 \\ -a, & a < 0 \end{array}\right.</code>		

122

Postscript/PDF graphics facilities

Applying coordinate transforms:

The `graphicx` package provides access to the geometric transform capabilities of the PostScript and PDF languages:

```
\scalebox{0.8}{\includegraphics{diagram.pdf}}
\includegraphics[height=60mm]{diagram2.pdf}
\resizebox{190mm}{60mm}{becomes 19 cm x 6 cm large}
\resizebox{190mm}{!}{this becomes 19 cm wide}
\rotatebox{180}{this is upside down!}
```

Changing colours:

The `color` package also uses Postscript/PDF special commands:

This text is `\textcolor{red}{printed in red}` if ...

This text is **printed in red** if you include `\usepackage{color}`.

Default: `\definecolor{red}{rgb}{1,0,0}`

124

Figures and references

Larger diagrams interfere with page breaking. They are best placed into a figure environment, such that \LaTeX can move them around. Example:

```
\begin{figure}
  \includegraphics[width=0.6\linewidth]{photo.jpg}
  \caption{This photograph shows the experimental setup.}
  \label{fig:expsetup}
\end{figure}
```

The automatically assigned figure number can be quoted as in:

```
See also Figure~\ref{fig:expsetup}
(page~\pageref{fig:expsetup}).
```

The `\label{...}` command can also be used after `\section{...}`, `\subsection{...}`, etc. and inside `\begin{equation} ... \end{equation}` to assign symbolic names to section and equation numbers, which can then be resolved via `\ref{...}` or `\pageref{...}`.

No need to manually renumber figures, sections, or equations!

125

Build tools for \LaTeX

To make sure `\label` references and tables of contents use the correct numbers, it may be necessary to call `latex` twice. It will output “Rerun to get cross-references right” in this case.

The following implicit Makefile rule takes care of this:

```
.DELETE_ON_ERROR:
%.pdf %.aux %.idx: %.tex
    pdflatex $<
    while grep 'Rerun to get ' $*.log ; do pdflatex $< ; done
```

An alternative is the “`latexmk`” tool, which automatically determines dependencies (e.g. from `\includegraphics`) and recompiles \LaTeX documents where file modification timestamps indicate that this is necessary.

126

Graphics editor xfig

- ▶ Its `*.fig` files have a simple plain-text format that can be edited manually, script generated, and leads to useful diffs.
- ▶ Can export `*.eps` or `*.pdf` files
- ▶ Can also produce figures in which \LaTeX is used to fill in all the text. This provides math mode, macros, symbols, references, fonts that match the main text, etc.

Ask `xfig` to export a `*.pstex + *.pstex_t` file pair. The `*.pstex` file lacks the text parts of the figure. The `*.pstex_t` file contains \LaTeX commands that first load the `*.pstex` image, and then add all the text in the figure. Select the “special text” mode in `xfig` to enable \LaTeX metacharacters. Use `\include{*.pstex_t}` to add such a figure in your document. (PDF equivalent: `*.pdftex + *.pdftex_t`)

- ▶ Command-line export tool (e.g., for Makefile): `fig2dev`

```
%.eps: %.fig
    fig2dev -L eps $< $@
%.pstex %.pstex_t: %.fig
    fig2dev -L pstex_t -p $*.pstex $< $*.pstex_t
    fig2dev -L pstex $< $*.pstex
%.pdftex %.pdftex_t: %.fig
    fig2dev -L pdftex_t -p $*.pdftex $< $*.pdftex_t
    fig2dev -L pdftex $< $*.pdftex
```

Other graphics tools: TikZ, pnmtools, Inkscape, MATLAB, R, gnuplot, Python+matplotlib

127

Conclusions

- ▶ Unix is a powerful and highly productive platform for experienced users.
- ▶ This short course could only give you a quick overview to get you started with exploring advanced Unix facilities.
- ▶ Please try out all the tools mentioned here and consult the “`man`” and “`info`” online documentation.
- ▶ You’ll find on
<http://www.cl.cam.ac.uk/teaching/current/UnixTools/>
easy to print versions of the `bash`, `make` and `perl` documentation, links to further resources, and hints for installing Linux on your PC.

★ ★ Good luck and lots of fun with your projects ★ ★

128