

Any large, useful software system ...

- will undergo many changes in its lifetime
- multiple programmers making changes
- who may work on the code concurrently and independently

The process of code change needs to be managed so that

- changes produce "consistent" versions of the system
- many programmers can easily work simultaneously
- we can roll back to earlier version if needed
- documentation of when, who, & why changes made
- multiple versions of system can be distributed, tested, merged

Consider the following simple scenario:

- a software system contains a source code file `x.c`
- system is worked on by several teams of programmers
- Ann in Sydney adds a new feature in her copy of `x.c`
- Bob in Singapore fixes a bug in his copy of `x.c`

Ultimately, we need to ensure that

- all changes are properly recorded (when, who, why)
- both the new feature and the bug fix are in the next release
- if we later find bugs in old release, they can be fixed

# Human Version Control

Manual solution to above problems:

- Ann requests/receives a copy of `x.c` v2.1 from Manager
- Bob requests/receives a copy of `x.c` v2.1 from Manager
- Ann sends her new version of `x.c` back to Manager
- Bob sends his new version of `x.c` back to Manager
- Manager ensures that all changes are incorporated
  - might need to ask Ann and Bob for help if changes conflict
- Manager sets up the merged version as v2.2 of `x.c`
- Col requests/receives a copy of `x.c` v2.2 from Manager

Problem: eventually Manager not available 24x7 & as system scales will be overwhelmed with requests.

# Version Control Systems

*Version control systems* aim to solve the above problems.

Version control systems (VCSs) are also called ...

- revision control systems
- source (code) control systems
- (source) code management systems

There are various approaches to solving the problems, leading to different families of version control systems. While VCSs could be used for all kinds of documents, we focus on their use for managing source code files.

# Version Control Systems

A *version control system* allows software developers to:

- share development work on a system
- recreate old versions of a system when needed
- identify the current versions of source code files
- restrict who is allowed to modify each source code file

This allows change to be managed/controlled in a systematic way. VCSs also try to minimise resource use in maintaining multiple versions.

# Labelling Versions

How to name an important version? (unique identifier)

Common approach: file name + version "number" (e.g. Perl 5.8.1)

No "standard" for *a.b.c* version numbers, but typically:

- *a* is major version number (changes when functionality/API changes)
- *b* is minor version number (changes when internals change)
- *c* is patch level (changes after each set of bug fixes are added)

Examples: Oracle 7.3.2, 8.0.5, 8.1.5, ...

# Labelling Versions

How to store multiple versions (e.g. v3.2.3 and v3.2.4)?

We *could* simply store one complete file for each version.

Alternative approach:

- store complete file for version 3.2.3
- store differences between 3.2.3 and 3.2.4

A *delta* is the set of differences between two successive versions of a file.

Many VCSs store a combination of complete versions and deltas for each file.

# Labelling Versions

Creating version  $N$  of file  $F$  ( $F_N$ ) from a collection of

- complete copies of  $F$  whose versions  $< N$
- deltas for all versions in between complete copies

is achieved via:

```
get list of complete copies of F
choose highest complete version V << N
f = copy of F{V}
foreach delta between V .. N {
    f = f + delta
}
# f == F{N} (i.e. version N of F)
```

Programs like patch can apply deltas.



# Delta Bandwidth Efficiency

An example of why deltas are useful:

- Google Chrome for Windows upgrades in the background almost daily
- Google Chrome is 10MB
- bsdiff delta = 700KB
- google custom delta (Courgette) = 80KB
- 200 full upgrades/year = 2GB/year
- 200 bsdiff upgrades/year = 140Mb/year
- 200 Courgette upgrades/year = 16Mb/year

# Unix VCS - Generation 1 (Unix)

1970's ... *SCCS* (source code control system)

- first version control system
- centralized VCS - single central repository
- introduced idea of multiple versions via delta's
- single user model: lock - modify - unlock
- only one user working on a file at a time

1980's ... *RCS* (revision control system)

- similar functionality to *SCCS*  
(essentially a clean open-source re-write of *SCCS*)
- centralized VCS - single central repository
- single user model: lock - modify - unlock
- only one user working on a file at a time
- still available and in use

# Unix VCS - Generation 2 (Unix)

1990 ... CVS (concurrent version system)

- centralized VCS - single central repository
- locked check-out replaced by copy-modify-merge model
- users can work simultaneously and later merge changes
- allows remote development essential for open source projects
- web-accessible interface promoted wide-dist projects
- poor handling of file metadata, renames, links

Early 2000's ... *Subversion* (svn)

- depicted as "CVS done right"
- many cvs weakness fixed
- solid, well documented, widely used system
- but essentially the same model as CVS
- centralized VCS - single central repository
- svn is suitable for assignments/small-medium projects
- easier to understand than distributed VCSs, well supported
- but Andrew recommends git

Early 2000s... *Bitkeeper*

- distributed VCS - multiple repositories, no "master"
- every user has their own repository
- written by Larry McVoy
- Commercial system but allowed limited use for Linux kernel until dispute over licensing issues
- Linus Torvalds + others then wrote GIT open source distributed VCS
- Other open source distributed VCS's appeared, e.g: bazaar (Canonical!), darcs (Haskell!), Mercurial

- distributed VCS - multiple repositories, no "master"
- every user has their own repository
- created by Linux Torvalds for Linux kernel
- external revisions imported as new branches
- flexible handling of branching
- various auto-merging algorithms
- Andrew recommends you use git unless good reason not to
- not better than competitors but better supported/more widely used (e.g. github/bitbucket)
- at first stick with a small subset of commands
- substantial time investment to learn to use Git's full power

Many VCSs use the notion of a *repository*

- store all versions of all objects (files) managed by VCS
- may be single file, directory tree, database,...
- possibly accessed by filesystem, http, ssh or custom protocol
- possibly structured as a collection of *projects*

# Git Repository

Git uses the sub-directory `.git` to store the repository.

Inside `.git` there are:

- **blobs** file contents identified by SHA-1 hash
- **tree objects** links blobs to info about directories, link, permissions (limited)
- **commit objects** links trees objects with info about parents, time, log message
- Create repository **git init**
- Copy existing repository **git clone**

# Tracking a Project with Git

- Project must be in single directory tree.
- Usually don't want to track all files in directory tree
- Don't track binaries, derived files, temporary files, large static files
- Use **.gitignore** files to indicate files never want to track
- Use **git add *file*** to indicate you want to track *file*
- Careful: **git add *directory*** will every file in *file* and sub-directories



- A git commit is a snapshot of all the files in the project.
- Can return the project to this state using **git checkout**
- Beware if you accidentally add a file with confidential info to git - need to remove it from all commits.
- **git add** copies file to staging area for next commit
- **git commit -a** if you want commit current versions of all files being tracked
- commits have parent commit(s), most have 1 parent
- merge produce commit with 2 parents, first commit has no parent
- merge commits labelled with SHA-1 hash

# Git Branch

- Git branch is pointer to a commit
- Allows you to name a series of commits
- Provides convenient tracking of versions of parallel version of projects
- New features can be developed in a branch and eventually merged with other branches
- Default branch is called master.
- **HEAD** is a reference to the last commit in the current branch
- **git branch** *name* creates branch *name*
- **git checkout** *name* changes all project files to their version on this branch

# Git Merge

- merges branches
- git mergetool shows conflicts
- configure your own mergetool - many choices kdiff3, meld, p4merge
- kaleidoscope popular on OSX

- **git push** *repository-name branch* adds commits from your *branch* to remote repository *repository-name*
- Can set defaults, e.g. **git push -u origin master** then run **git push**
- **git remote** lets you give names to other repositories
- Note **git clone** sets *origin* to be name for cloned repository

# Git Fetch/Pull

- **git fetch** *repository-name branch* adds commits from *branch* in remote repository *repository-name*
- Usually **git pull** - combines fetch and merge

## Example: Git & conflicts

```
% git init /home/cs2041
Initialized empty Git repository in /home/cs2041/.git/
% git add main.c
% git commit main.c
Aborting commit due to empty commit message.
% git commit main.c -m initial
[master (root-commit) 8c7d287] initial
1 files changed, 1 insertions(+), 0 deletions(-)
create mode 100644 main.c
```

# Example: Git & conflicts

Suppose Fred does:

```
% cd /home/fred
% git clone /home/cs2041/.git /home/fred
Cloning into /home/fred...
done.
% echo >fred.c
% git commit -m 'created fred.c'
```

# Example: Git & conflicts

Suppose Jane does:

```
% cd /home/jane
% git clone /home/cs2041/.git /home/jane
Cloning into /home/jane...
done.
% echo '/* Jane Rules */' >>main.c
% git commit -a -m 'did some documentation'
[master 1eb8d32] did some documentation
1 files changed, 1 insertions(+), 0 deletions(-)
```



## Example: Git & conflicts

Fred can now get Jane's work like this:

```
% git pull /home/jane/.git
remote: Counting objects: 5, done.
remote: Total 3 (delta 0), reused 0 (delta 0)
Unpacking objects: 100% (3/3), done.
From /home/jane/.git
 * branch                HEAD          -> FETCH_HEAD
Merge made by recursive.
 main.c |      1 +
 1 files changed, 1 insertions(+), 0 deletions(-)
```

## Example: Git & conflicts

And Jane can now get Fred's work like this:

```
% git pull /home/fred/.git
remote: Counting objects: 7, done.
remote: Compressing objects: 100% (4/4), done.
remote: Total 5 (delta 0), reused 0 (delta 0)
Unpacking objects: 100% (5/5), done.
From /home/fred/.git
 * branch                HEAD          -> FETCH_HEAD
Updating 1eb8d32..63af286
Fast-forward
 fred.c |      1 +
 1 files changed, 1 insertions(+), 0 deletions(-)
 create mode 100644 fred.c
```

# Example: Git & conflicts

But if Fred does this:

```
% echo '// Fred Rules' >fred.c  
% git commit -a -m 'added documentation'
```

And Jane does this:

```
% echo '// Jane Rules' >fred.c  
% git commit -a -m 'inserted comments'
```

# Example: Git & conflicts

When Fred tries to get Jane's work:

```
% git pull /home/jane/.git
remote: Counting objects: 5, done.
remote: Compressing objects: 100% (2/2), done.
remote: Total 3 (delta 0), reused 0 (delta 0)
Unpacking objects: 100% (3/3), done.
From ../jane/
* branch                HEAD          -> FETCH_HEAD
Auto-merging fred.c
CONFLICT (content): Merge conflict in fred.c
Automatic merge failed; fix conflicts and then commit the r
% git mergetool
```

# Example: making Git Repository Public via Github

Github popular repo hosting site (see competitors e.g. bitbucket)  
Github free for small number of public repos  
Github and competitors also let you setup collaborators, wiki, web pages, issue tracking  
Web access to git repo e.g. <https://github.com/mirrors/linux>

# Example: making Git Repository Public via Github

Its a week after the 2041 assignment was due and you want to publish your code to the world.

Create github account - assume you choose *2041rocks* as your login

Create a repository - assume you choose *my\_code* for the repo name

Add your ssh key (`.ssh/id_rsa.pub`) to github (Account Settings - SSH Public Keys - Add another public key)

```
cd ~/cs2041/ass2
```

```
git remote add origin git@github.com:2041rocks/my_code.git
```

```
git push -u origin master
```

Now anyone anywhere can clone your repository by

```
git clone git@github.com:2041rocks/my_code.git
```

# Simple Version Control System

You can implement in a small amount of Perl the basic operations of a VCS. We'll use these utility functions:

```
sub read_file {
my ($file) = @_;
    open(my $f, '<', $file) or die "Can not read '$file': $!\n";
    return do {local $/; <$f>}
}

sub write_file {
my ($file, $contents) = @_;
    open my $f, '>>', $file or die "Can not write '$file': $!\n";
    print $f $contents;
}

sub read_dir {
my ($dir) = @_;
my @files = glob "$dir/*";
s/.*\/// for each @files;
return @files;
}
```

Git uses a directory tree in `.git` by default to store the repository. Our simple VCS (`sit.pl`) uses `.sit`, storing each version in a separate numbered subdirectory.



# Creating an empty repository:

- git init
- sit.pl init

Our simple VCS just needs to create a directory and initialize a file.

```
die "repository already exists\n" if -d $repository_dir;
mkdir $repository_dir or
    die "$0: can not create $repository_dir:$!";
write_file($commit_file, 0);
print "Created empty repository\n";
```

# Placing Files under Version Control

- `git add FILES`
- `git.pl add FILES`

Our simple VCS copies the files into a sub-directory of the repository (the stage area)

Git does also stages the files for commit.

```
-d $staging_area or mkdir $staging_area or  
die "$0: can not create $staging_area: $!\n";  
write_file("$staging_area/$_",read_file($_)) foreach @ARGV;  
my @staged_files = read_dir("$staging_area");  
print "Files now staged: @staged_files\n";
```

# Saving a Version

- `git commit -m LOG_MESSAGE`
- `sit.pl commit LOG_MESSAGE`

Git labels commits with a SHA-1 hash (160 bit 40 hex digits).  
Our simple VCS labels each commit with consecutive integer.  
It moves the files from its staging area into a sub-directory named with this number and writes the log message.

```
die "No files staged for commit\n" if !-d $staging_area;
my $last_commit = read_file($commit_file);
my $commit = $last_commit + 1;
write_file($commit_file, $commit);
$version_dir = "$repository_dir/$commit";
rename $staging_area, $version_dir;
write_file("$repository_dir/log_message.$commit", "@ARGV");
my @committed_files = read_dir("$repository_dir/$commit");
print "Committed as commit #$commit: @committed_files\n";
```

# Checking out a version

Common to want to examine a particular version of a file.

- `git show file [branch]`
- `sit.pl show file [version]`

Our simple VCS has to work backwards through versions from the requested version to find the appropriate version of the file.

```
my $file = shift @ARGV or die "$0: show must specify a file\n";
my $checkout_commit = shift @ARGV || read_file($commit_file);
# go through commits in reverse order finding last commit of e
foreach $commit (reverse 1..$checkout_commit) {
    my $path = "$repository_dir/$commit/$file";
    next if !-r $path;
    print STDERR "Commit #$commit: $file\n";
    print read_file($path);
    exit(0);
}
die "$0: no commit of $file\n";
```

# Reviewing commits

- git log [huge number of options]
- sit.pl log

Our simple VCS just prints the log message for each commit and which files it contained.

```
my $last_commit = read_file($commit_file);
foreach $commit (1..$last_commit) {
    print "Commit #$commit\n";
    print "Log message: ", read_file("$repository_dir/log_r
my @committed_files = read_dir("$repository_dir/$commit"
    print "Files: @committed_files\n";
}
```

# Checking Status of Files

- git status
- sit.pl log

Our simple VCS indicates files in the current directory not in the repository or which are different to the latest version in the repository.

```
my $last_commit = read_file($commit_file);
FILE: foreach $file (glob "*") {
    my $path = "$staging_area/$file";
    if (-r $path) {
        if (read_file($file) ne read_file($path)) {
            print "Modified since staged: $file\n";
        } else {
            print "Staged for commit: $file\n";
        }
    } else {
        foreach $commit (reverse 1..$last_commit) {
            my $path = "$repository_dir/$commit/$file";
            next if !-r $path;
            if (read_file($file) ne read_file($path)) {
                print "Modified and unstaged: $file\n";
            } else {
                print "Unmodified since committed: $file\n";
            }
            next FILE;
        }
        print "Untracked: $file\n";
    }
}
```

# Inspecting file differences

- git diff [many options]
- sit.pl unimplemented

# Updating your copy from repository

- git pull (or fetch)
- sit.pl unimplemented



# Using our simple Simple Version Control System

```
% sit.pl init
Created empty repository
% sit.pl add *. [ch]
Files staged are: world.c world.h graphics.c graphics.h main.c
% sit.pl commit 'initial import'
Committed as commit #1: world.c world.h graphics.c graphics.h main.c
% echo >>world.c
% echo >>graphics.c
% sit.pl add world.c
% echo >>world.c
% sit.pl status
Untracked: Makefile
Unstaged: graphics.c
Modified: world.c
% sit.pl add world.c graphics Makefile
Files staged are: world.c graphics.c Makefile
% sit.pl commit 'Minor changes'
Committed as commit #2: world.c graphics.c Makefile
% sit.pl log
Commit #1
Log message: initial import
Files: world.c world.h graphics.c graphics.h main.c
Commit #2
Log message: Minor changes
Files: world.c graphics.c Makefile
```