

Use Study in System and Data-Structure Design

distributed version-control system, apparently the most used currently.

Git stores snapshots (*versions*) of the files and directory structure of a project, keeping track of their relationships, authors, and commit messages.

Git is designed, in that there can be many copies of a given repository, allowing independent development, with machinery to transmit and receive versions between repositories.

Git is extremely fast (as these things go).

Major User-Level Features (I)

Git is a graph of versions or snapshots (called *commits*) of the project.

The structure reflects ancestry: which versions came from which.

Each commit contains

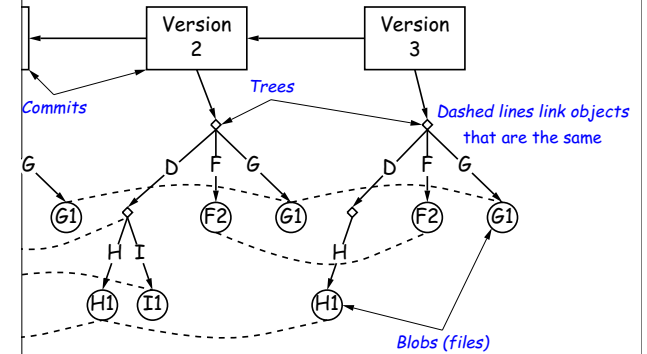
1. A pointer to a tree of files (like a Unix directory).

2. Information about who committed and when.

3. A pointer to the previous commit.

4. A pointer to the commit (or commits, if there was a merge) from which it was derived.

Commits, Trees, Files



Lecture #32

A Little History

Git was created by Linus Torvalds and others in the Linux community when they were dissatisfied with their previous, proprietary VCS (Bitkeeper) and wanted a free and open alternative.

The development effort seems to have taken about 2-3 months, leading to the 2.6.12 Linux kernel release in June, 2005.

The name, according to Wikipedia,

was chosen because Torvalds has quipped about the name Git, which is British slang meaning "unpleasant person". Torvalds said: "I'm a bit of a git, and I name all my projects after myself. 'Linux', now 'git'." The man page describes Git as "the distributed version control system."

Git is implemented as a collection of basic primitives (now called "plumbing") which are scripted to provide desired functionality.

Higher-level commands ("porcelain") built on top of these to provide a convenient user interface.

Conceptual Structure

Git objects consist of four types of *object*:

1. *Tree*: A directory structure of files.

2. *Blob*: A file's contents.

3. *Commit*: A snapshot of the repository at a specific point in time, containing references to trees and additional information (author, date, log message).

4. *Tag*: A reference to a commit, used to identify releases, other important versions, or other useful information. (Won't mention further today).

Major User-Level Features (II)

has a name that uniquely identifies it to all versions.
can transmit collections of versions to each other.

by a commit from repository A to repository B requires
transmission of those objects (files or directory trees)
not yet have (allowing speedy updating of repositories).

maintain named *branches*, which are simply identifiers
commits that are updated to keep track of the most
its in various lines of development.

branches are essentially named pointers to particular commits.
branches in that they are not usually changed.

The Pointer Problem

it are files. How should we represent pointers between

able to *transmit* objects from one repository to another
nt contents. How do you transmit the pointers?

transfer those objects that are missing in the target
how do we know which those are?

counter in each repository to give each object there a
But how can that work consistently for two independent

How A Broken Idea Can Work

to use a hash function that is so unlikely to have a
we can ignore that possibility.

Secure Hash Functions have relevant property.

tion, f , is designed to withstand cryptanalytic attacks.
, should have

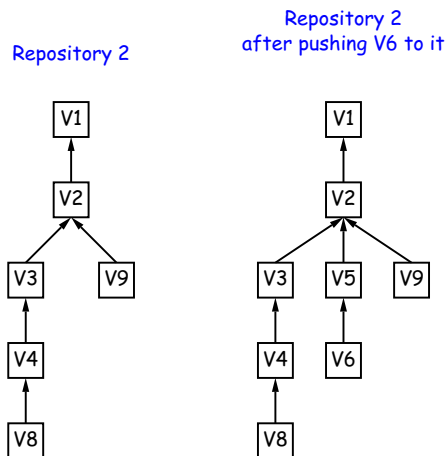
Pre-image resistance: given $h = f(m)$, should be computationally
to find such a message m .

Second-pre-image resistance: given message m_1 , should be infeasible
 $m_2 \neq m_1$ such that $f(m_1) = f(m_2)$.

Collision resistance: should be difficult to find *any* two messages
such that $f(m_1) = f(m_2)$.

properties, scheme of using hash of contents as name is
likely to fail, even when system is used maliciously.

Version Histories in Two Repositories



Internals

repository is contained in a directory.

may either be *bare* (just a collection of objects and
r may be included as part of a working directory.

the repository is stored in various *objects* corresponding
their "leaf" content), trees, and commits.

e, data in files is *compressed*.

git gc *prune* the objects from time to time to save additional

Content-Addressable File System

one way of naming objects that is universal.

names, then, as pointers.

Which objects don't you have?" problem in an obvious

, what is invariant about an object, regardless of repository,
hashes.

: the contents as the name for obvious reasons.

hash of the contents as the address.

at doesn't work!

1: Use it anyway!!

Low-Level Blob Management

out the hashcode that *Git* uses for the blob containing
eg. java with the command

```
git hash-object something.java
```

tells you that the file would have hash code

```
0d159f1550b0b5e102f7e06867cc44782
```

ally `git add` this file, its compressed contents will be
e file

```
objects/19/2a0ca0d159f1550b0b5e102f7e06867cc44782
```

ook at them (uncompressed) with

```
git file -p 192a0ca0d159f1550b0b5e102f7e06867cc44782
```

SHA1

1 (Secure Hash Function 1).

and with this using the `hashlib` module in Python3.

mes in *Git* are therefore 160-bit hash codes of contents,

commit in the shared CS61B repository could be fetched
with

```
git checkout 3b30599cc43f4616eb626f8fa4fb2d0610d97963
```