

git

an introduction



How Communication Channels Shape and Challenge a Participatory Culture in Software Development

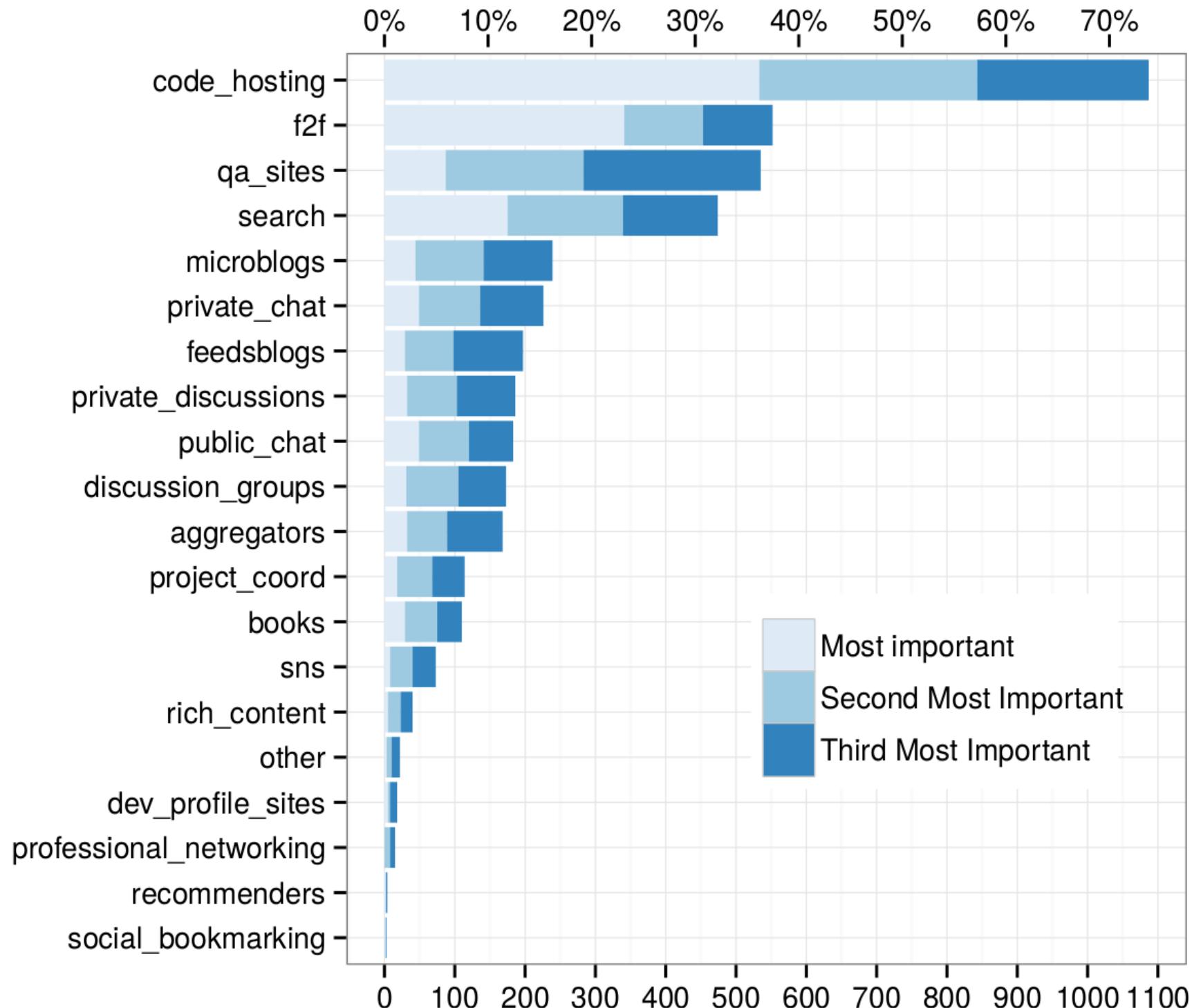
Margaret-Anne Storey, Leif Singer, Fernando Figueira Filho, Alexey Zagalsky, and Daniel M. German

Abstract—Software developers use many different communication tools and channels in their work. The diversity of these tools has dramatically increased over the past decade, giving rise to a wide range of socially-enabled communication channels and social media that developers use to support their activities. A participatory culture of software development is emerging in which developers want to engage with, learn from, and co-create software with other developers. However, the interplay of these channels, as well as the opportunities and challenges they may create when used together, are not yet well understood.

In this paper we report on a large-scale survey conducted with 1,448 GitHub users. We describe which channels these developers find essential to their work, and gain an understanding of the challenges they face using them. Our findings lay the empirical foundation for providing recommendations to developers and tool designers on how to use and improve tools for developers.

Index Terms—Participatory Culture, Communication, Social Media, CSCW, Software Engineering.

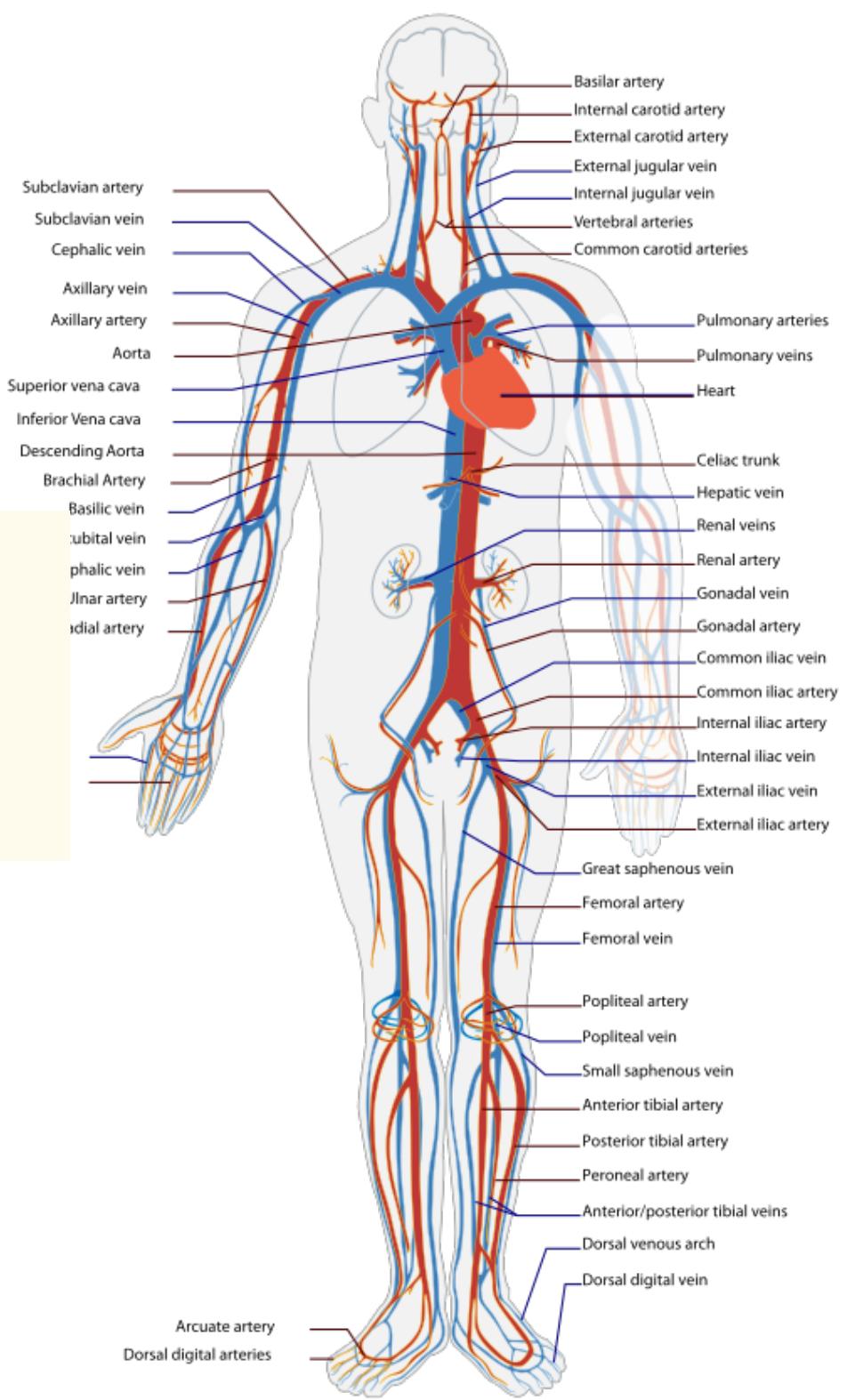




**How is distributed version control
affecting software development?**

VERSION CONTROL

a developer's best friend



1997

*"The most important book about technology today,
with implications that go far beyond programming."*
—Guy Kawasaki

THE CATHEDRAL & THE BAZAAR

MUSINGS ON LINUX AND OPEN SOURCE
BY AN ACCIDENTAL REVOLUTIONARY



ERIC S. RAYMOND

WITH A FOREWORD BY BOB YOUNG, CHAIRMAN & CEO OF RED HAT, INC.

The *Bazaar* model

Eric S Raymond credits Linus Torvalds as the inventor of the Bazaar model:

“the Linux community seemed to resemble a great babbling bazaar of differing agendas and approaches out of which a coherent and stable system could seemingly emerge”

Surprisingly, Linux **was not using version control** software at the time

**Linus/Linux has an interesting
relationship with version control**

**“Tarballs and patches is a much
superior version control system
than CVS is”**

**“I see Subversion as the most
pointless project ever started”**

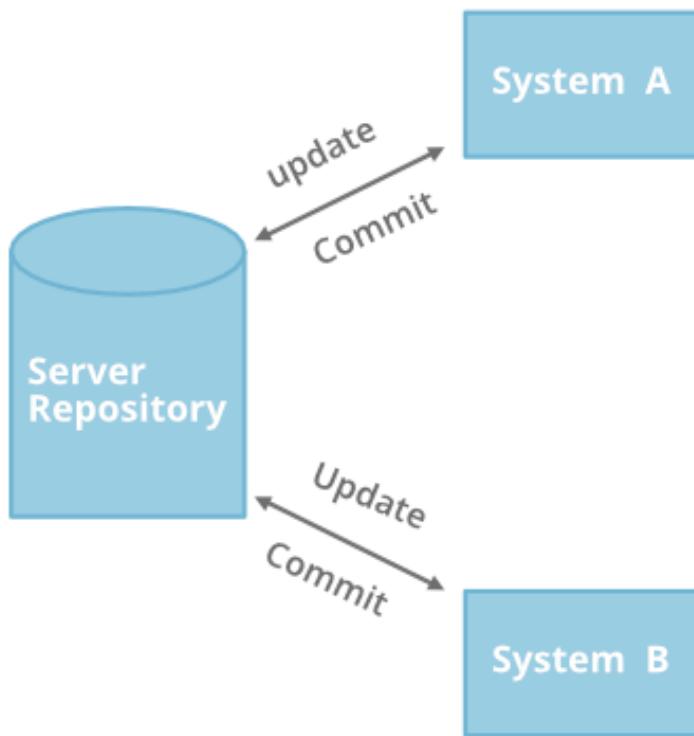
Centralized Version Control

- Centralized version control relies on a central, publicly accessible repository
 - Any developer who wants to participate, needs to have a user account

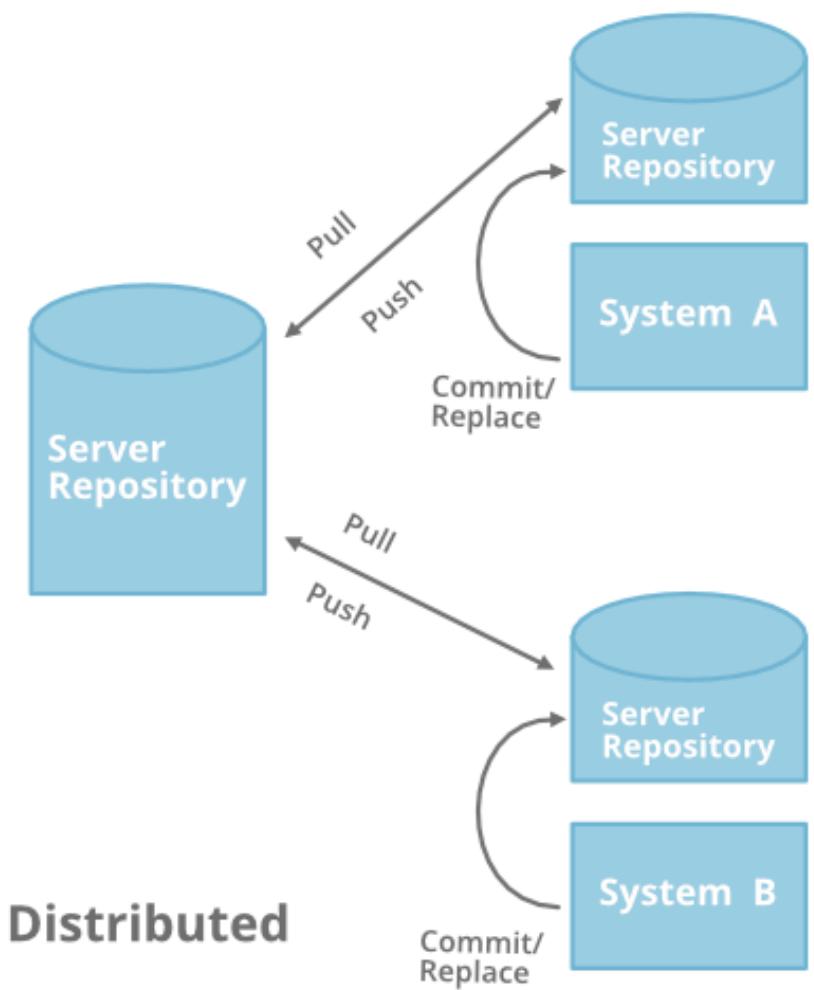
It means that before you start contributing, somebody needs to trust you

Distributed Version Control

- There are multiple repositories, all equally important
 - No need to “**register**” to contribute
 - Easy to start creating a feature
 - Relies on effective merging mechanisms



Centralized



Distributed

Implicit vs Explicit Branching

- In C-VCs
 - One must create a branch explicitly to work independently
 - Needs continuous access to the central repo for operations
- In D-VCs
 - Any repo is a branch of another repo
 - All operations are local

**“[source control management]
was uninteresting and yet
profoundly important.”**

“I hated all SCMs with a passion”

‘The big thing about distributed source control is that it makes one of the main issues with SCM’s go away – the politics around “who can make changes.”’

"So git was basically designed and written for my requirements, and it shows."

“I’m just happy that it made it
so **easy** to start a new project.
Project hosting used to be painful,
and with git and GitHub it’s just so
trivial to do a random small project.”

Bitkeeper

- Linus adopted Bitkeeper as its first version control system in 2002
 - It was not open source
 - Used until 2005, when the license is revoked
 - This triggered the development of **git**
 - Bitkeeper influenced the development of git

Git

- Linus design goals:
 - “What would CVS would never do?”
 - Support distributed bitkeeper-style workflow
 - Protect against corruption and data loss
 - Be fast

Collaboration

- Isolation
 - Every developer commits without disturbing others
 - You do not have to trust anybody else
 - **Everybody has their own branch**
 - **Everybody has commit access**

Concurrent development!

But git need at least one public repo

- In order to collaborate
- In 2008, GitHub is created to fill this gap

git Internals and data model

The blob

- At its core, git is a storage system
- The blob is its basic unit of storage
- Every blob is addressable by its id (a hash)
- The simplest blob is a file

5b1d3..	
blob	size
<pre>#ifndef REVISION_H #define REVISION_H #include "parse-options.h" #define SEEN (1u<<0) #define UNINTERESTING (1u #define TREESAME (1u<<2)</pre>	

The tree

- A directory/folder is stored as a list of blobs
 - Which can be other directories

c36d4..

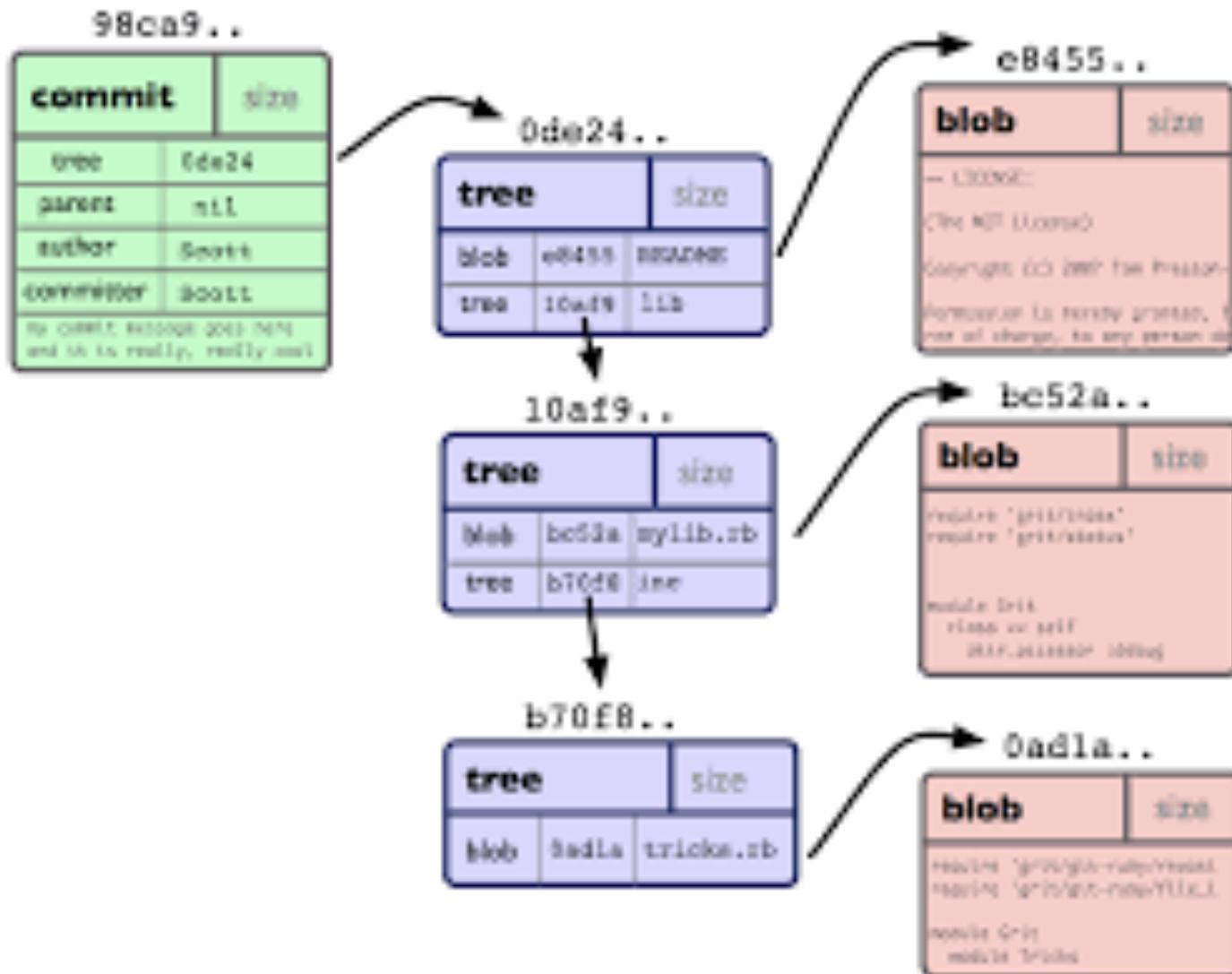
tree	size
blob 5b1d3	README
tree 03e78	lib
tree cdc8b	test
blob cba0a	test.rb
blob 911e7	xdiff

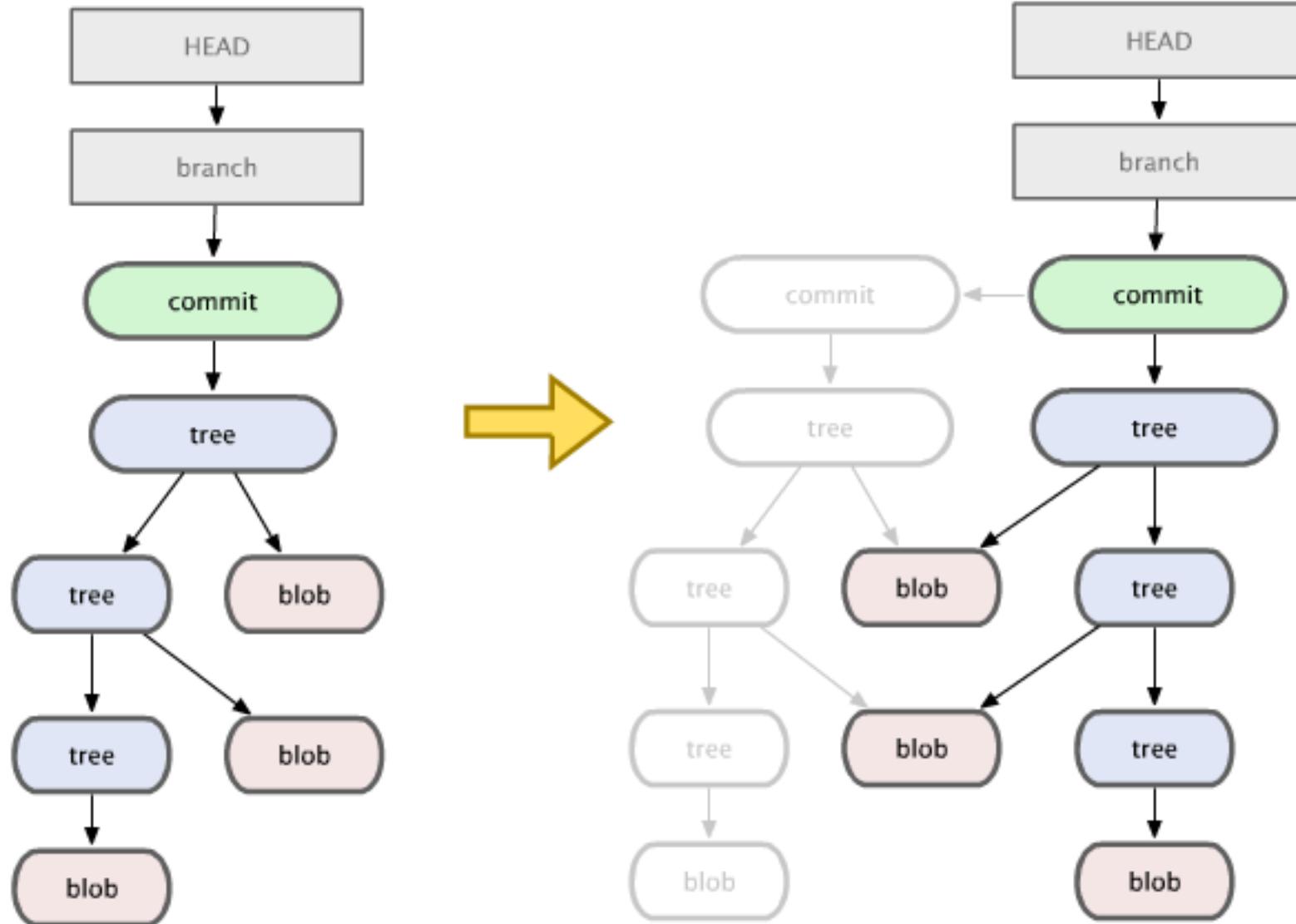
The commit

- Store:
 - the current tree (and its blobs)
 - The parent(s) commit(s)
- Along with some metadata
 - Author/Date
 - Committer/Date
 - Message

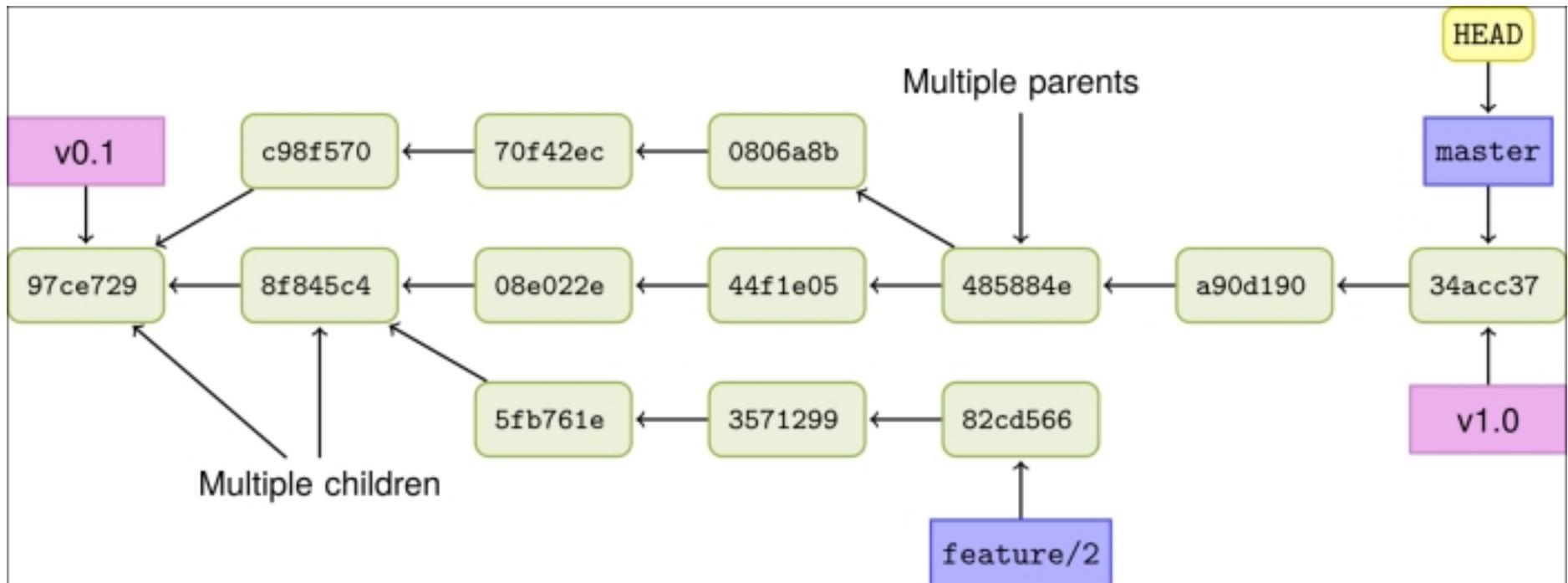
ae668..

commit	size
tree	c4ec5
parent	a149e
author	Scott
committer	Scott
my commit message goes here and it is really, really cool	





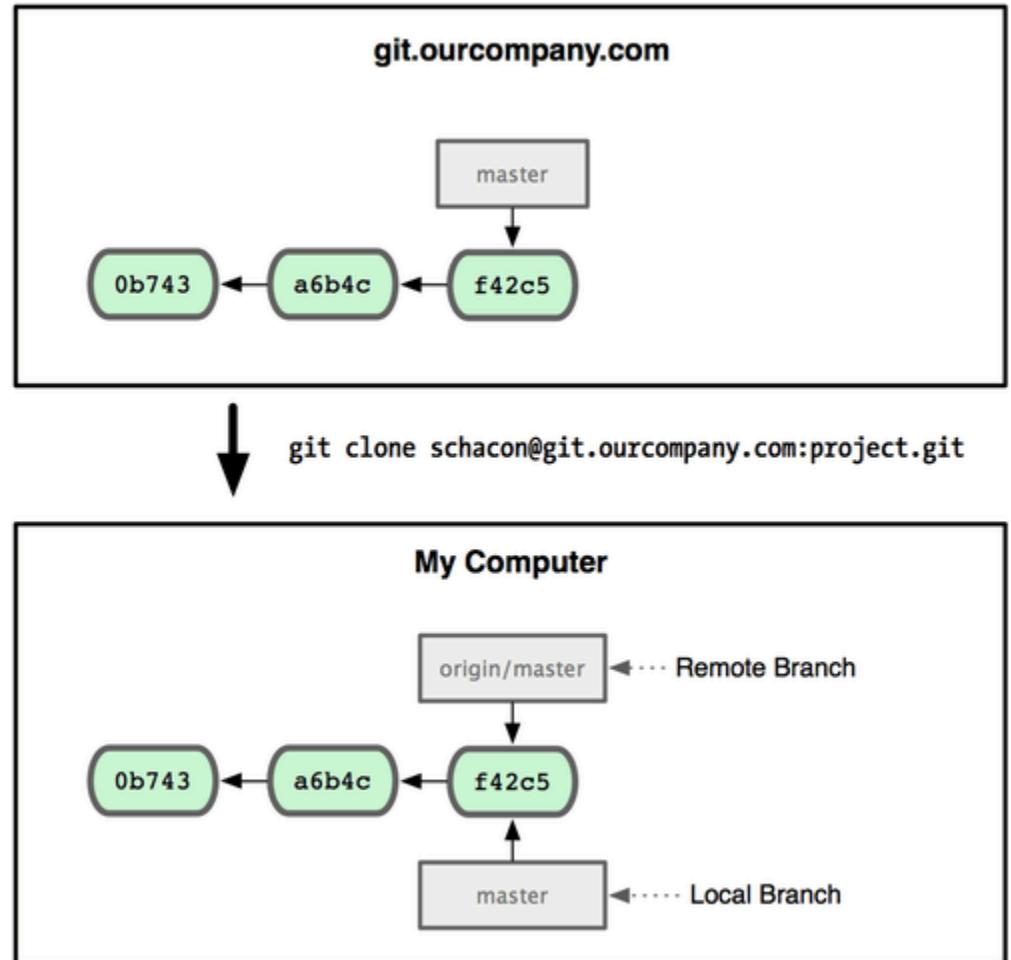
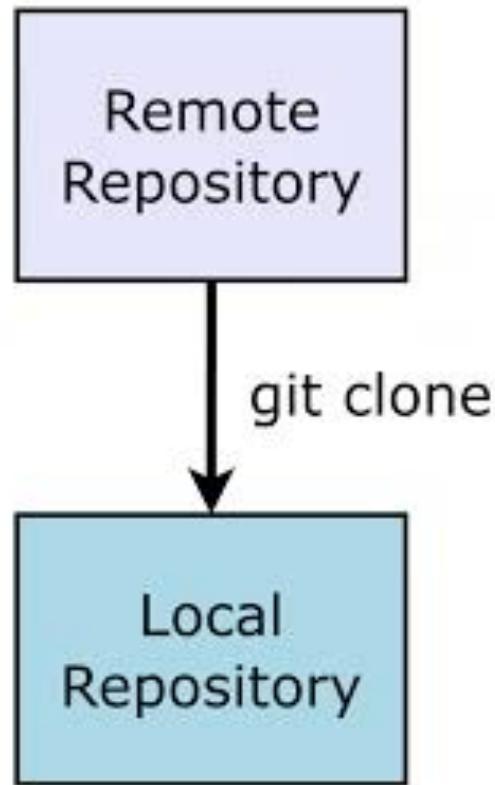
The DAG



Tags identify branches and specific commits

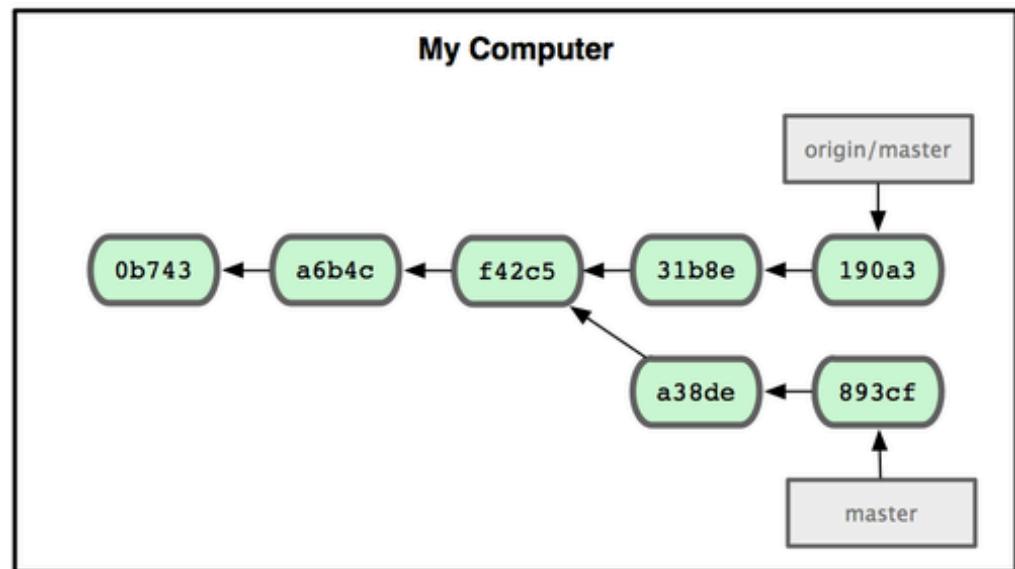
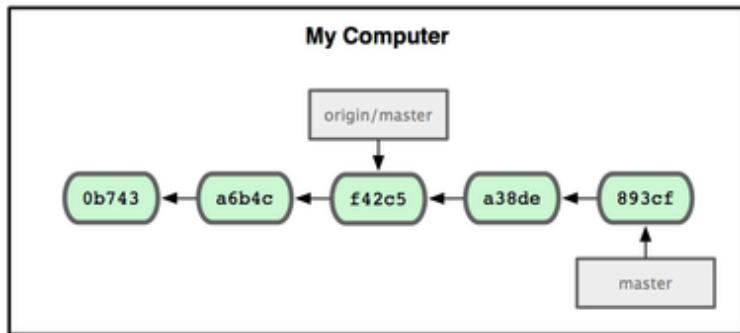
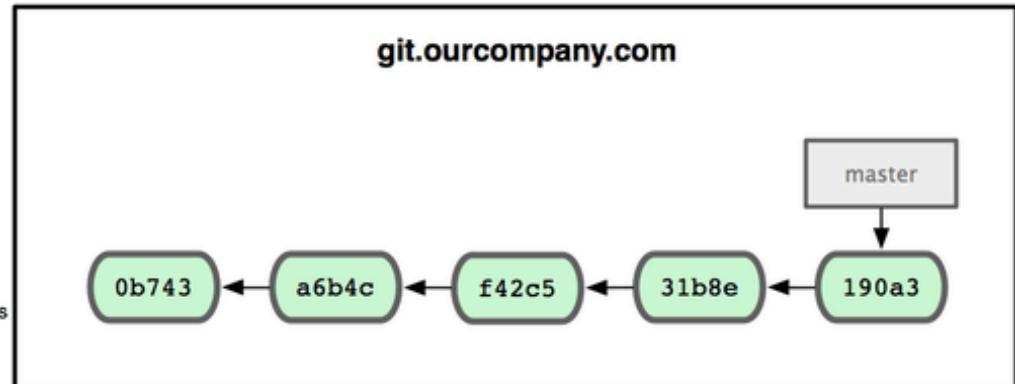
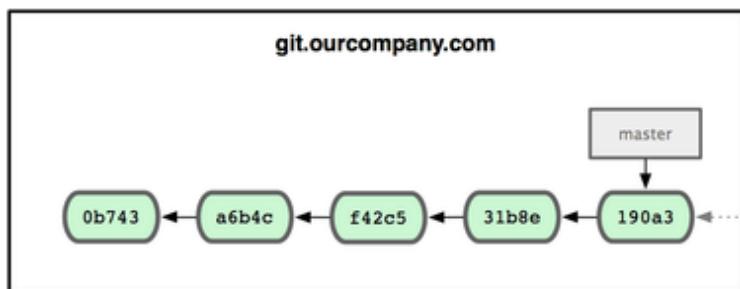
The operations

clone

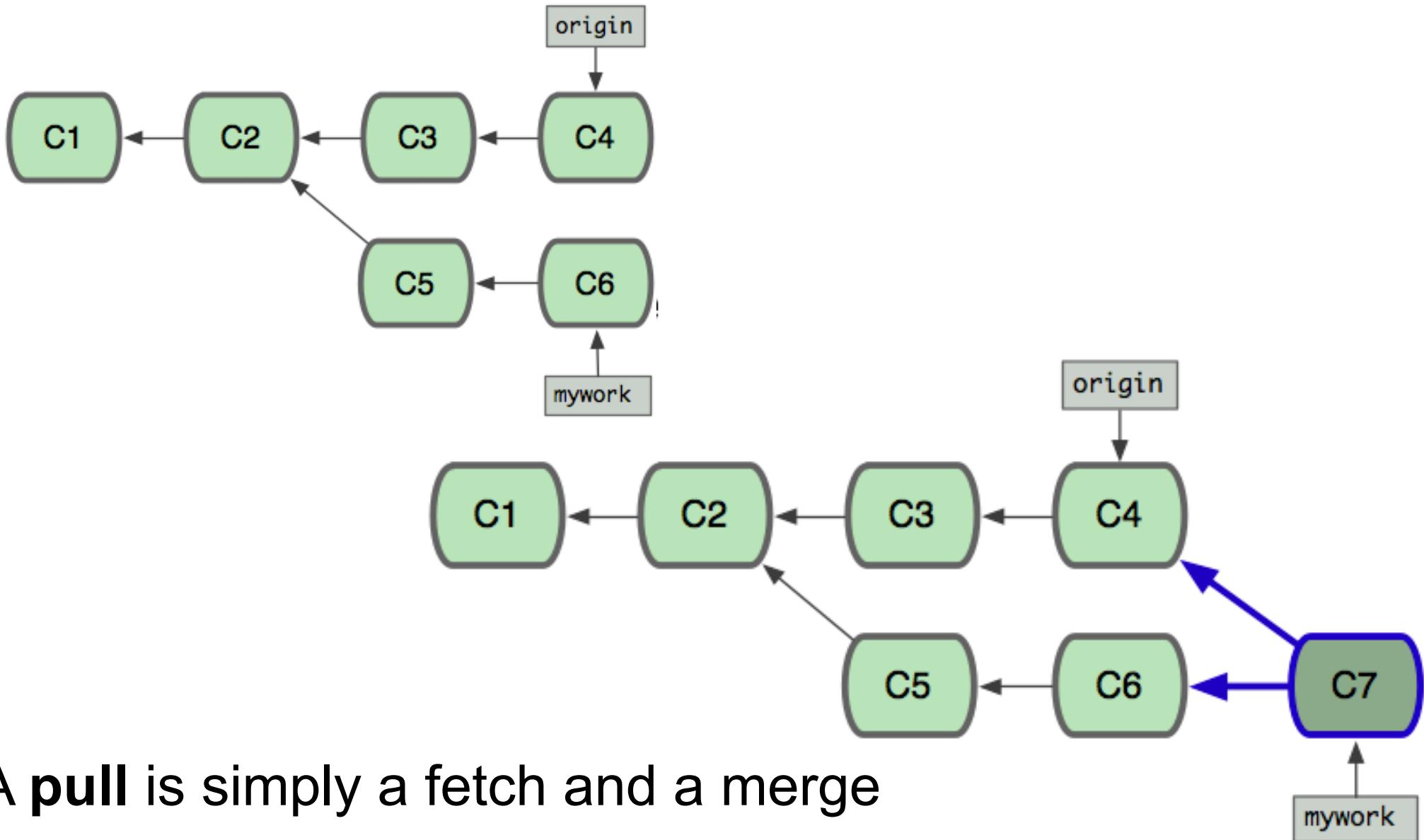


(git) clone is different from a (github) fork

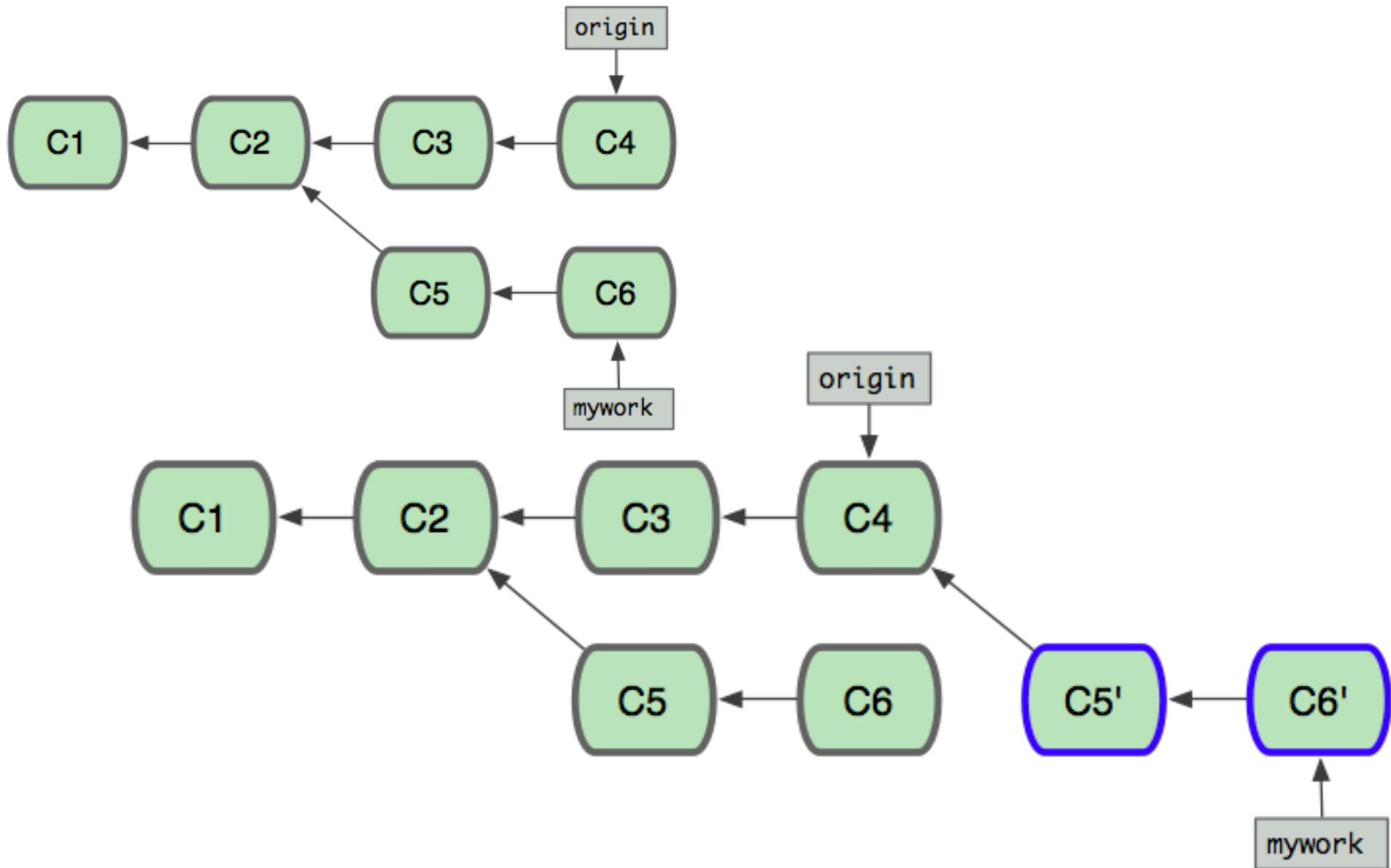
Fetch



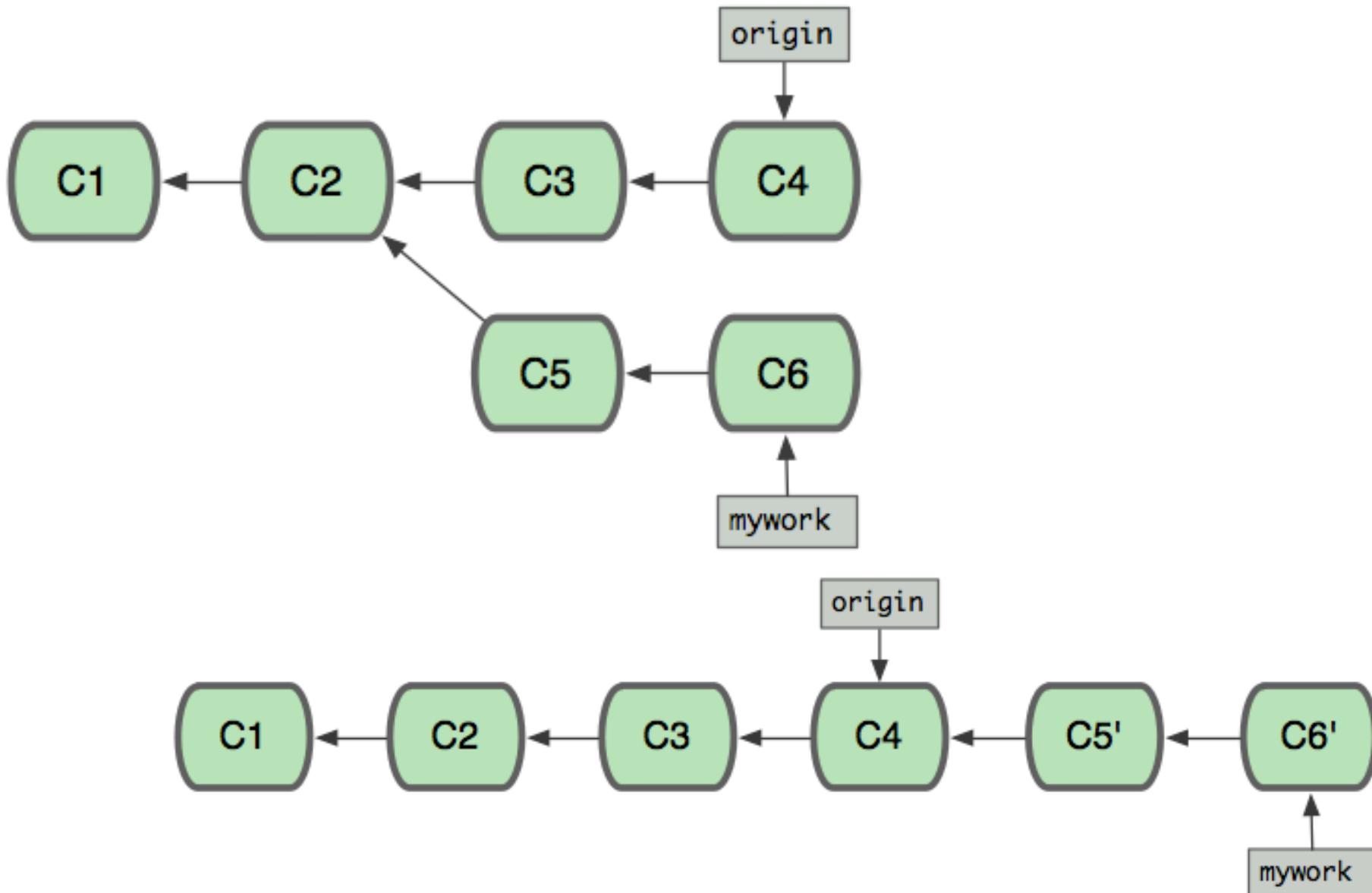
Merge



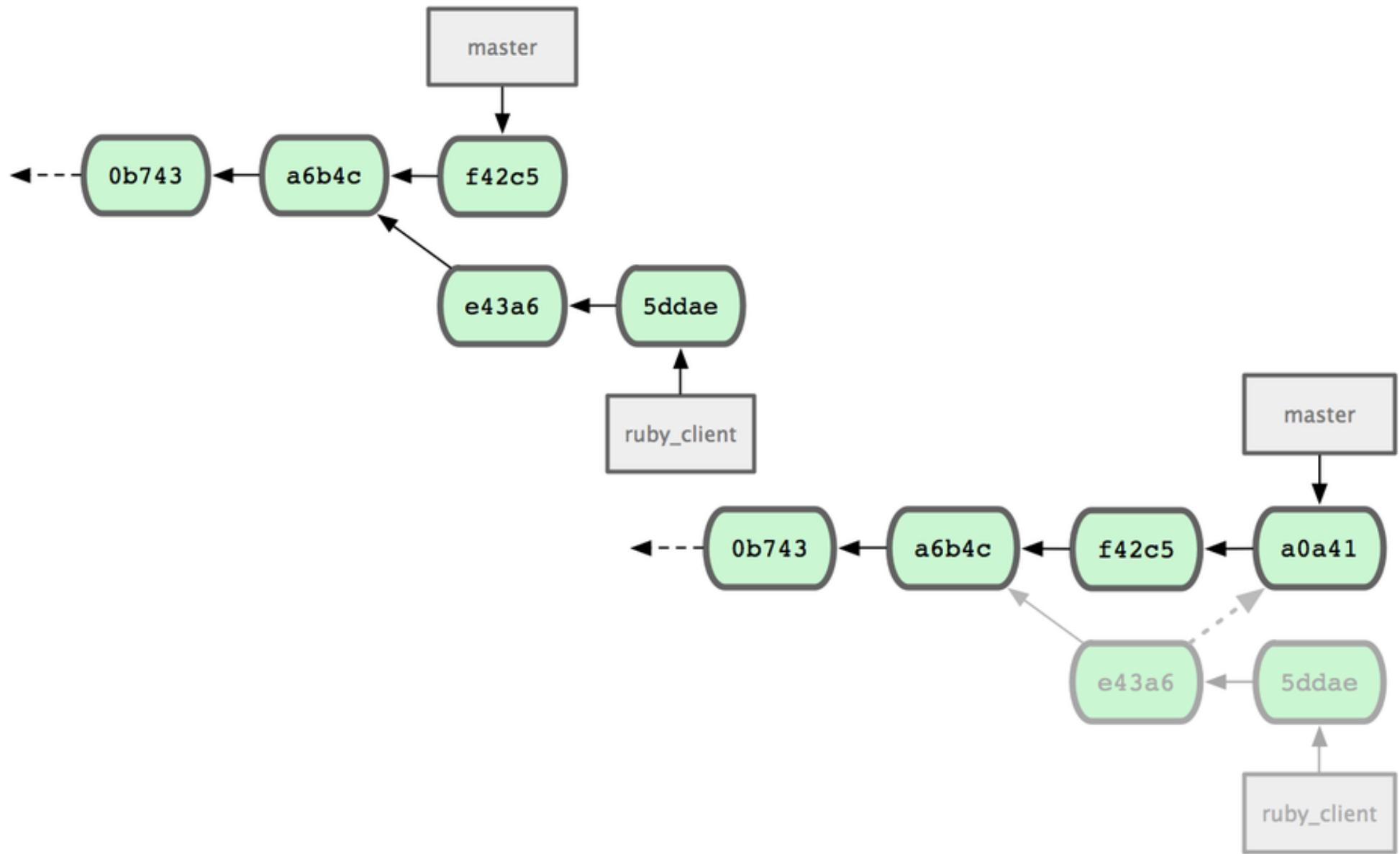
The rebase



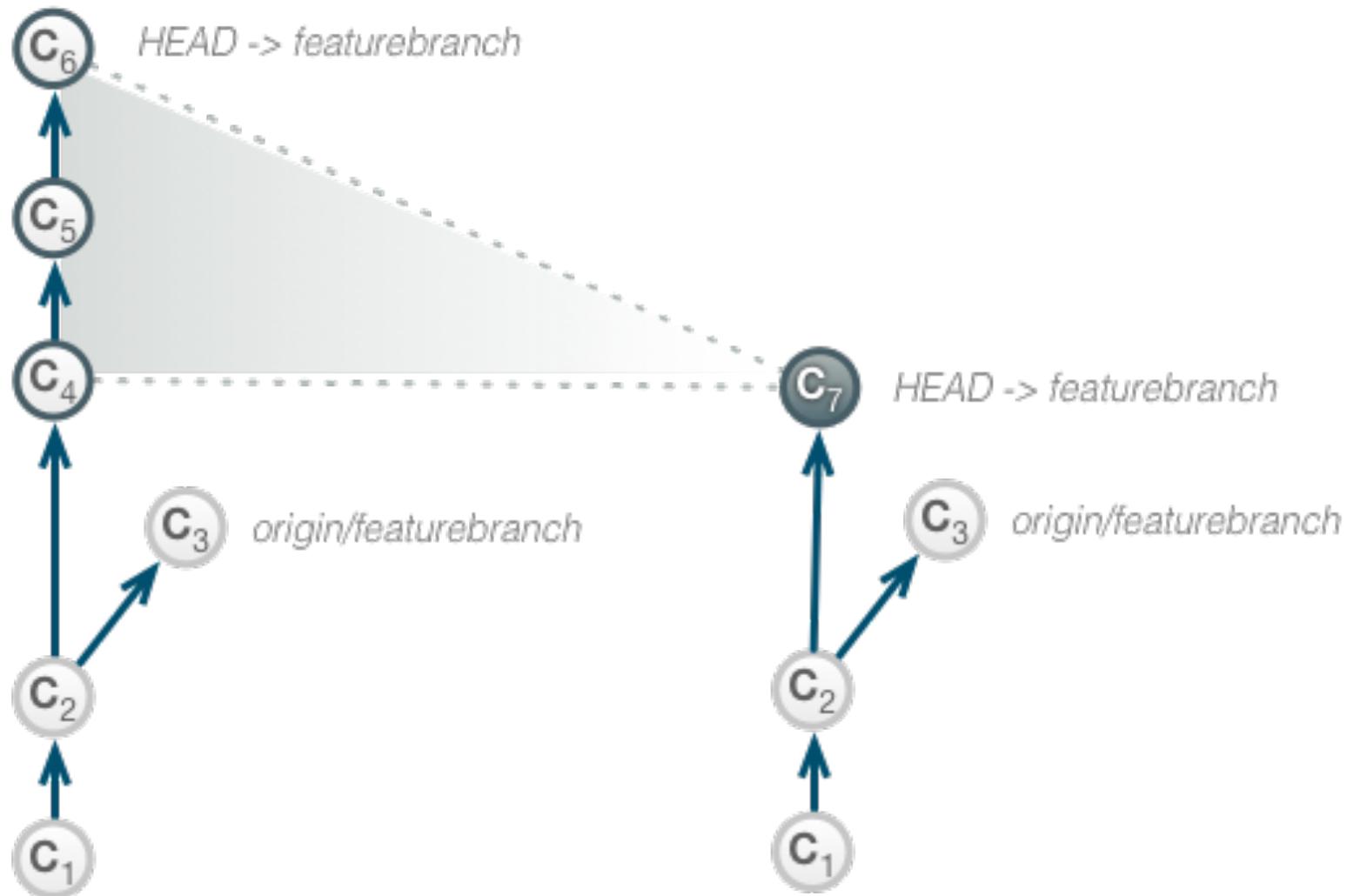
A fast-forward merge



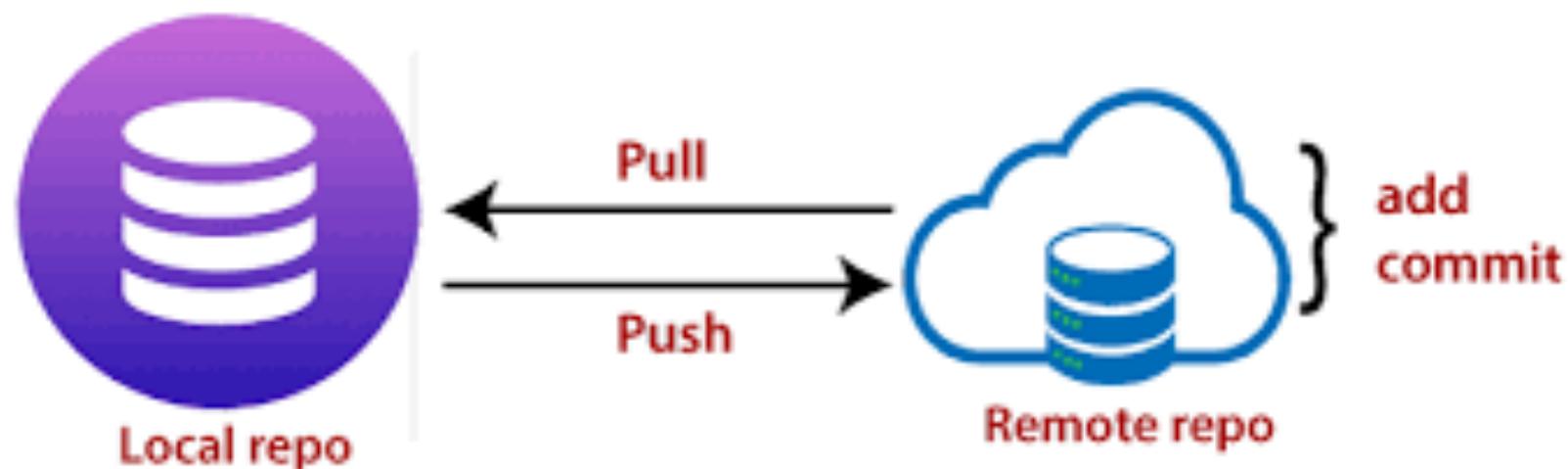
The cherry-pick



The squash



Push



How to propagate commits

- Using email
 - Create one or more emails with the patches
 - Recipient can take those emails and add their commits to their repo
- Push
 - Requires write access to the destination repo (usually public)
- Pull
 - Requires read access to the source repo (usually public)

Pull request

- A developer pushes changes to a public repo
- Then asks another developer to pull those changes to his/her local repository
- If accepted, the changes are then pushed to the public repository of the recipient
 - GitHub removes the need to do the pull request locally
 - It can be done from one public repo to another

Committer/Author?

- The author is the one that “writes” the patch
- The committer is the one that “commits” it to the current repo
- Most of the time, it is the
- Both author and committer
 - Name + email

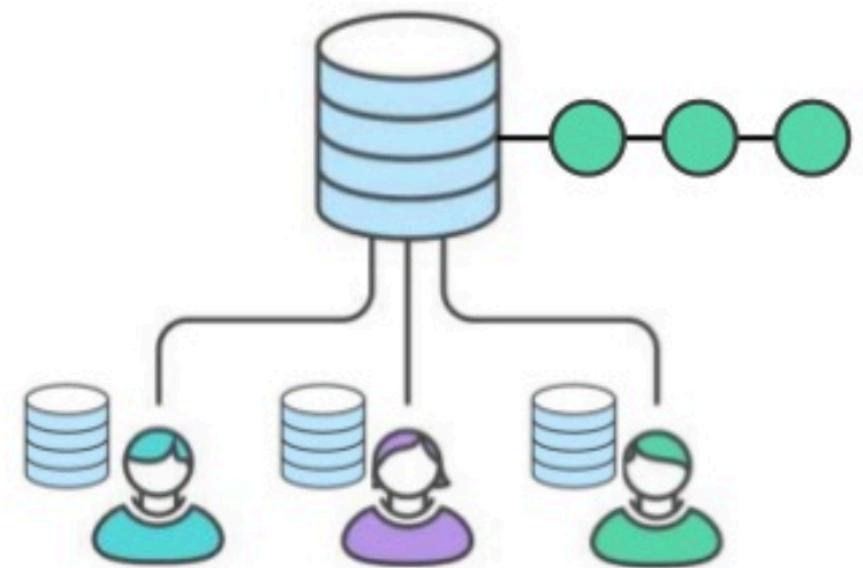
Several operations “recommit” a change

- Rebase, Squash and cherry-pick
 - The resulting commit keeps its author, but the committer is the one that does the change
 - If the commit has the same date for author and commit
 - Then the commit has never been rebased/squashed/cherry picked

The workflows

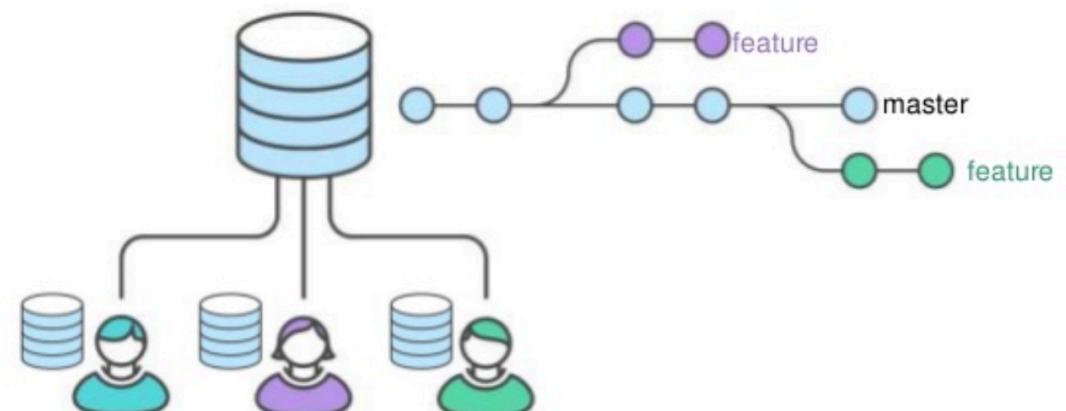
Centralized

- Each developer has a local repository
- But there is only one public repository
 - All developers fetch/pull/push to the same repo
 - No branch
 - No code review



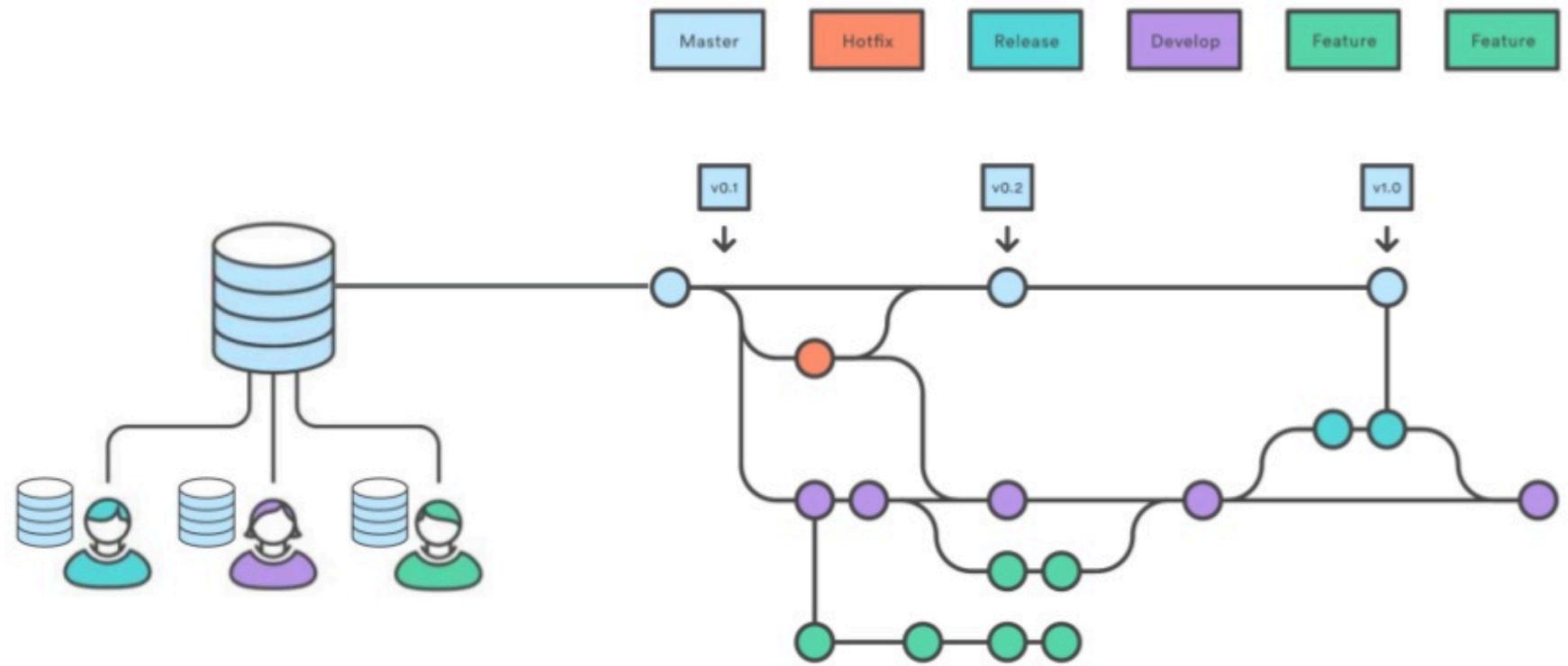
Feature branch

- Variant of centralized (one public repository)
 - Development is done in different branches
 - A branch is merged when the feature is completed
 - At that point, it disappears



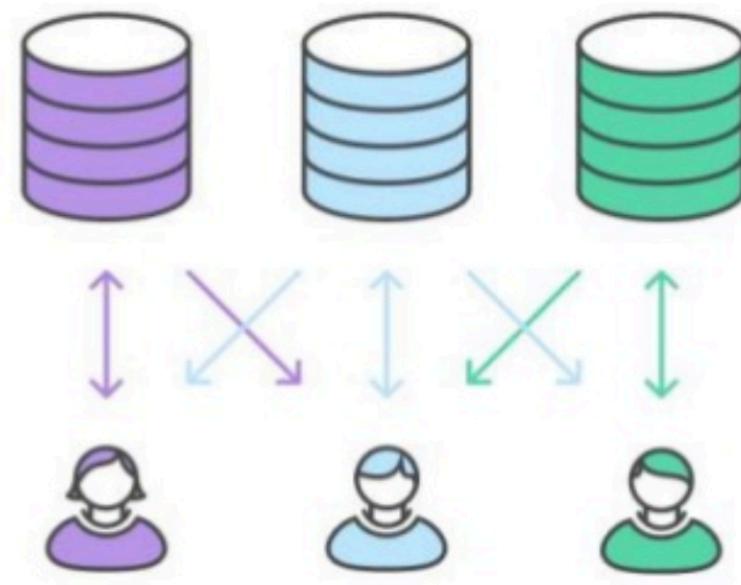
Git-flow

- Variant of centralized (one public repository)
 - More complex branching/cherry picking
 - Branches used for versions



Fork

- Multiple public repos
 - One repo is designed as the official
 - Code flows across repos





The Promises and Perils of Mining Git

Christian Bird*, Peter C. Rigby†, Earl T. Barr*, David J. Hamilton*, Daniel M. German†, Prem Devanbu*

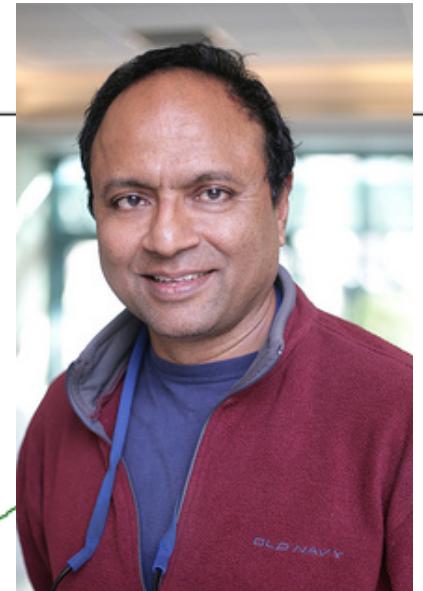
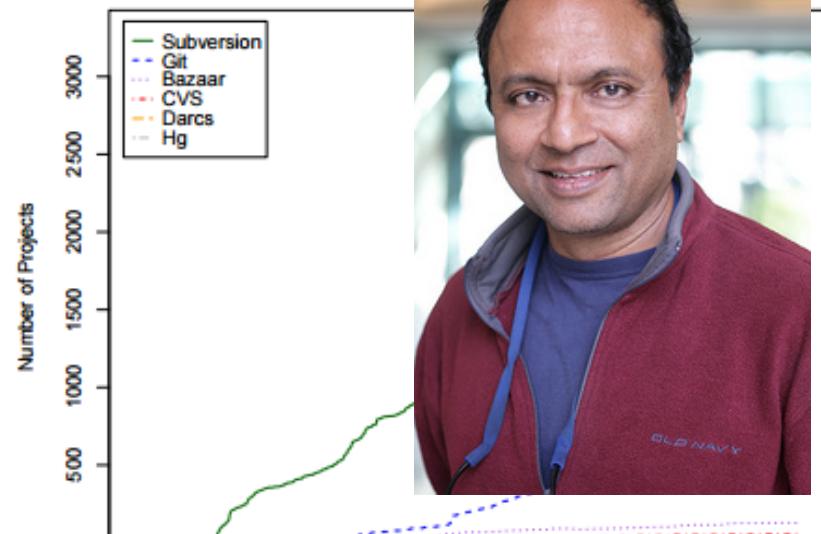
*University of California, Davis, USA

†University of Victoria, Canada

{bird, barr, hamiltod, devanbu}@cs.ucdavis.edu {pcr, dmger}@cs.uvic.ca

Abstract

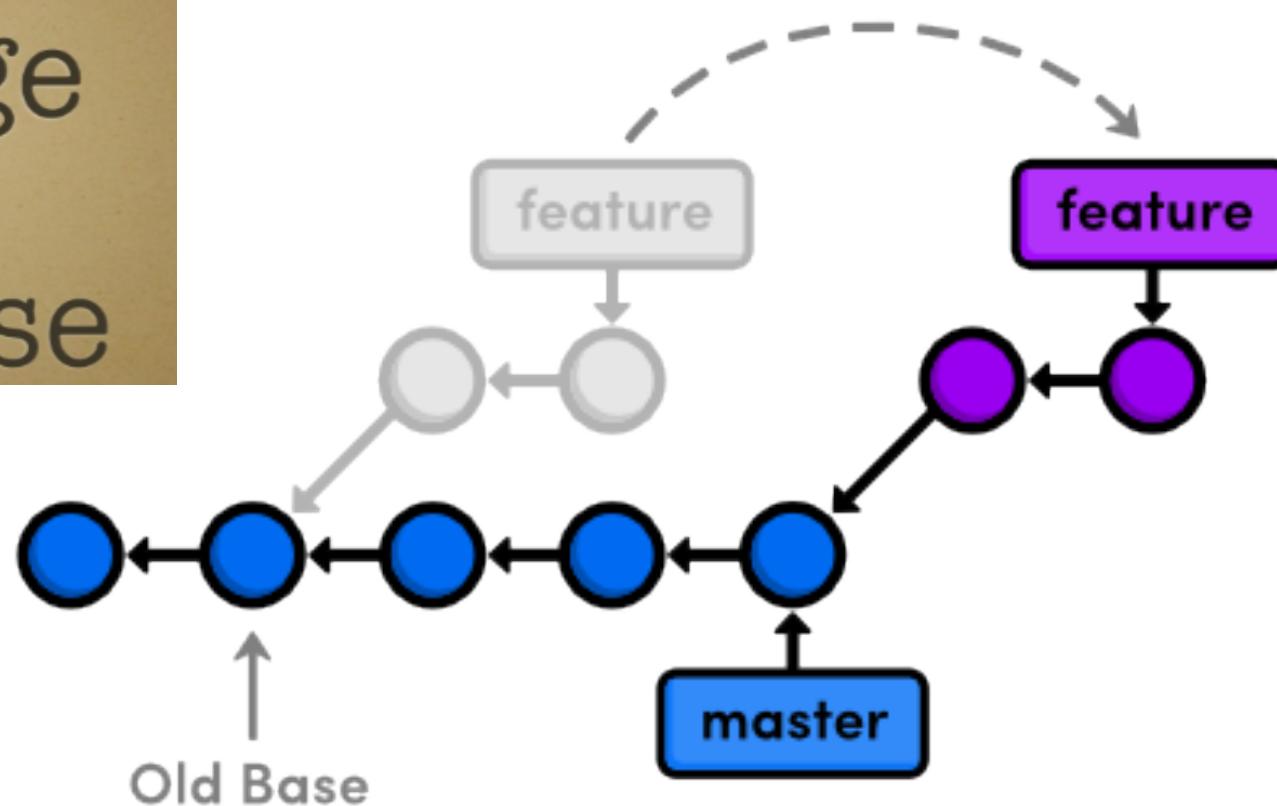
We are now witnessing the rapid growth of decentralized source code management (DSCM) systems, in which every developer has her own repository. DSCMs facilitate a style of collaboration in which work output can flow sideways (and privately) between collaborators, rather than always end down (and publicly) via a central repository. Decentralization comes with both the **promise** of new data and the **peril** of its misinterpretation. We focus on git, a very popular DSCM used in high-profile projects. Decentralization, and other features of git, such as automatically recorded contributor attribution, lead to richer content histories, giving



Git's “history”

The Rebase

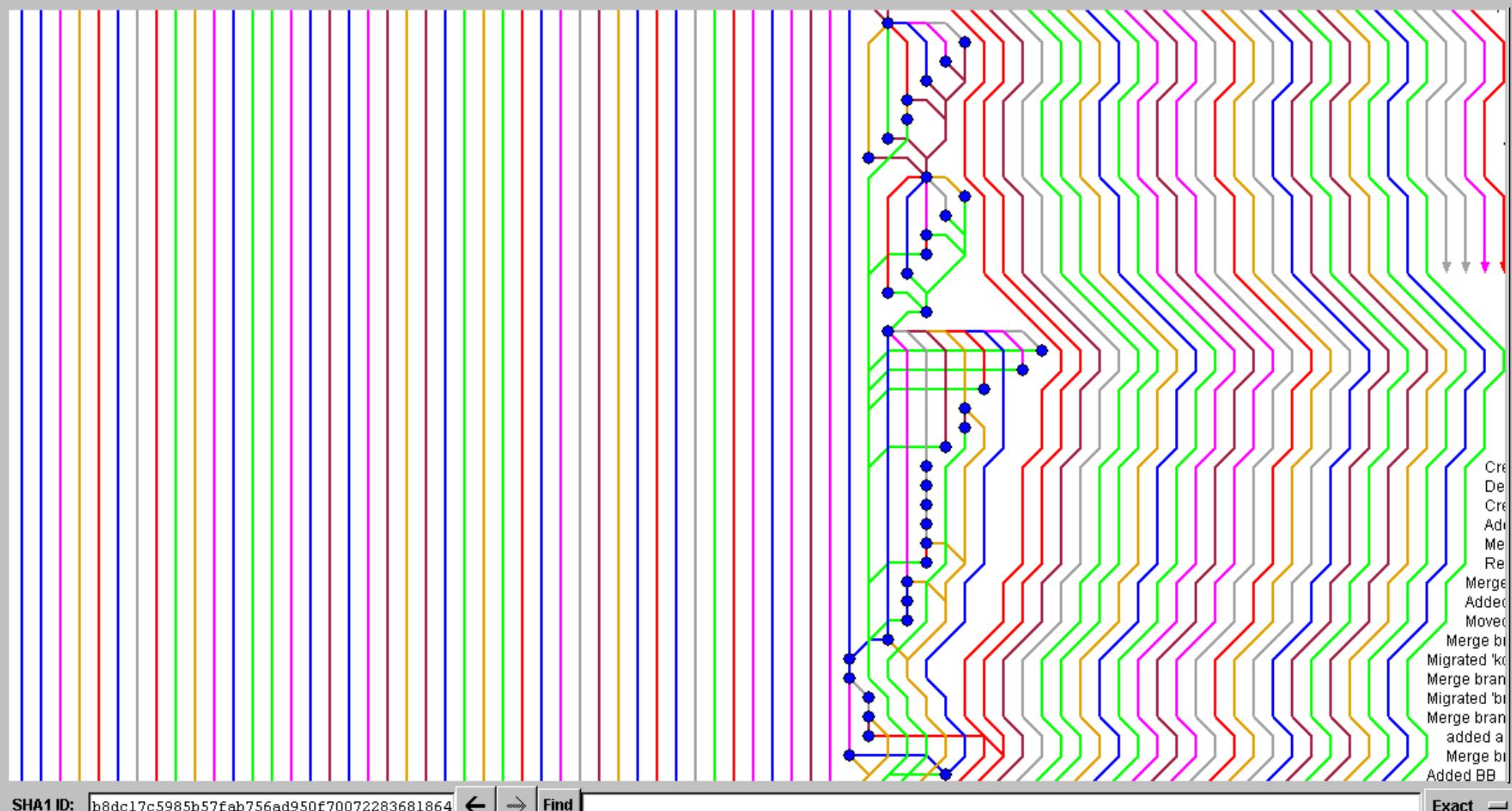
Merge
vs
Rebase



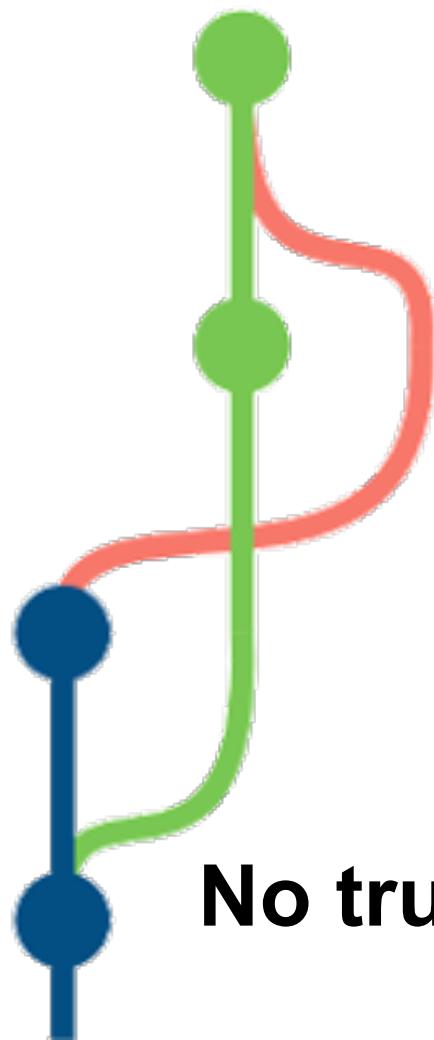
Where is the trunk?



There is no trunk



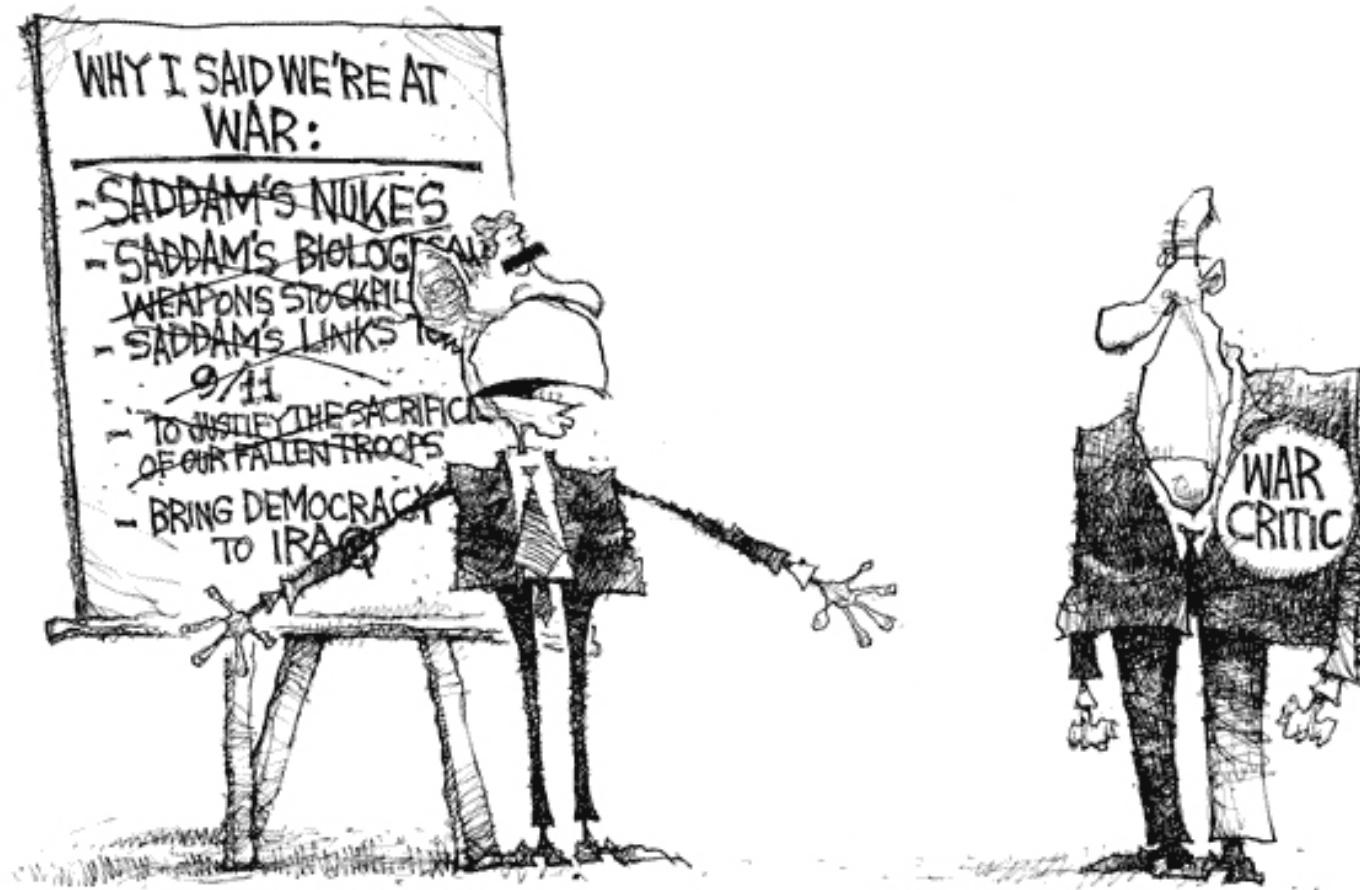
The Foxtrot



No true concept of “master”
branch

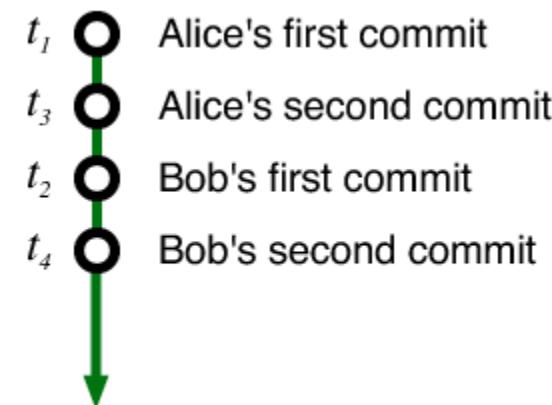
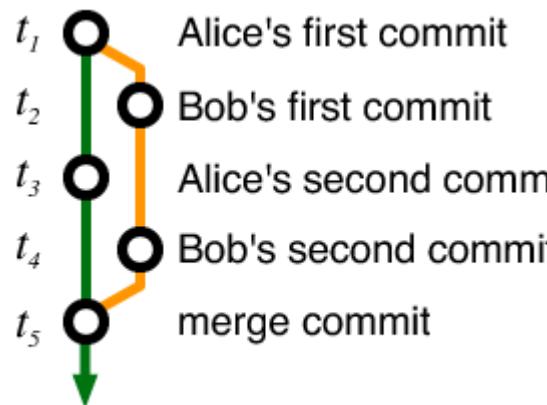
History? What is history?

THE JOURNAL NEWS
©2005 MAT DAVIES
#13



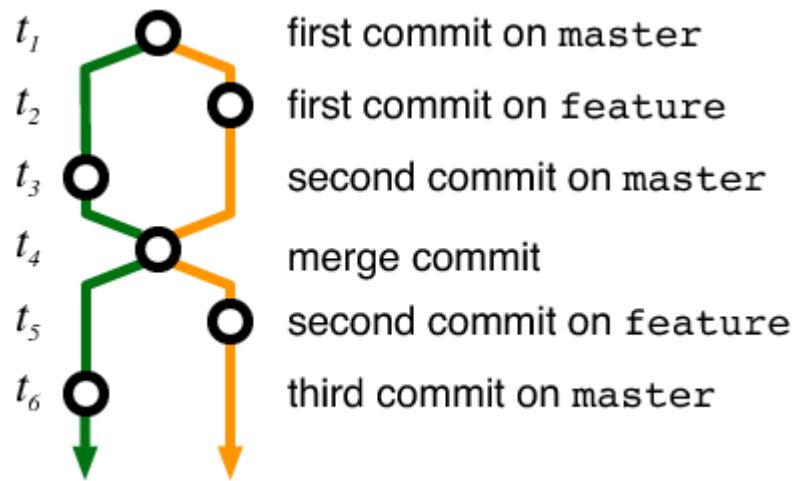
"...I WON'T LET YOU REWRITE HISTORY..."

Rebasing



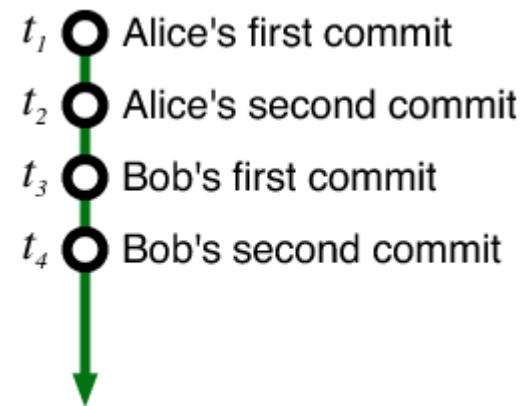
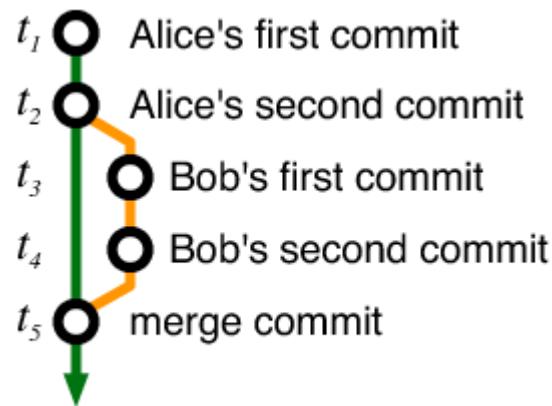
Makes for cleaner history!

In which repo?



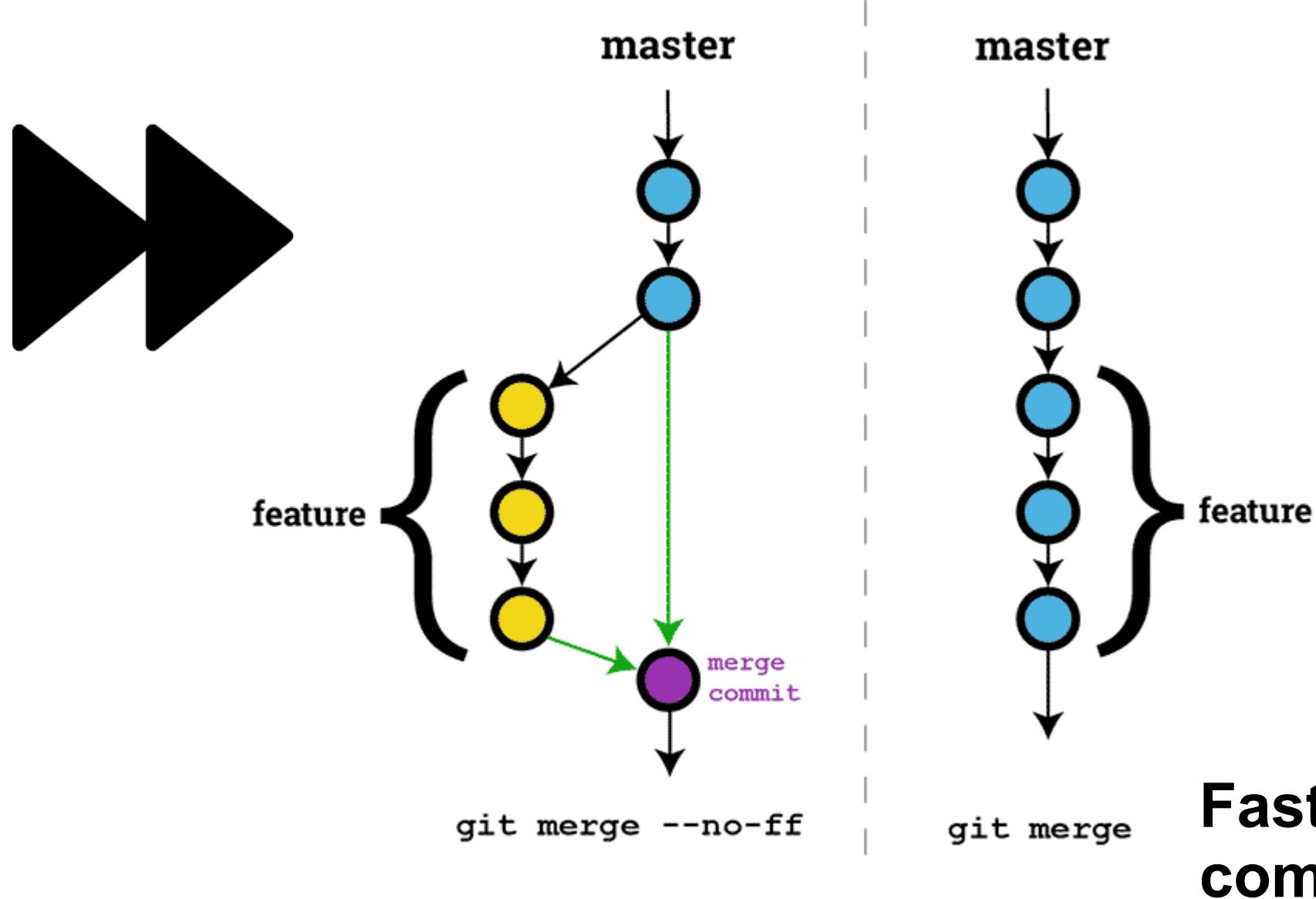
Nobody knows!

Was there a merge? I didn't see it!



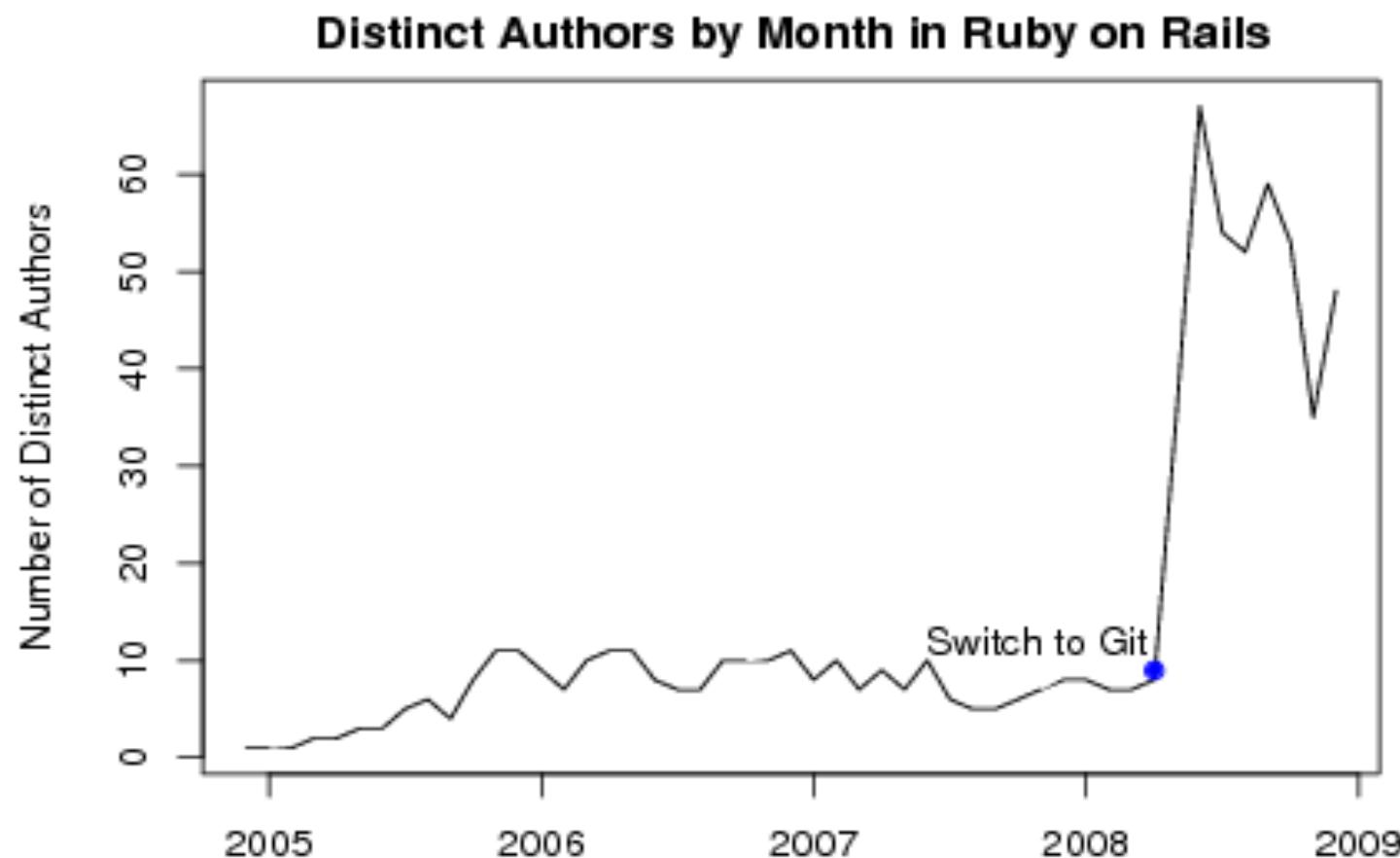
Fast-forward commits

Was there a merge? I didn't see it



Metadata

- Committer is not the same as the author



No centralized logging



Example: postgres

- It does **not** have a single merge
 - It is used as CVS



Linux and git

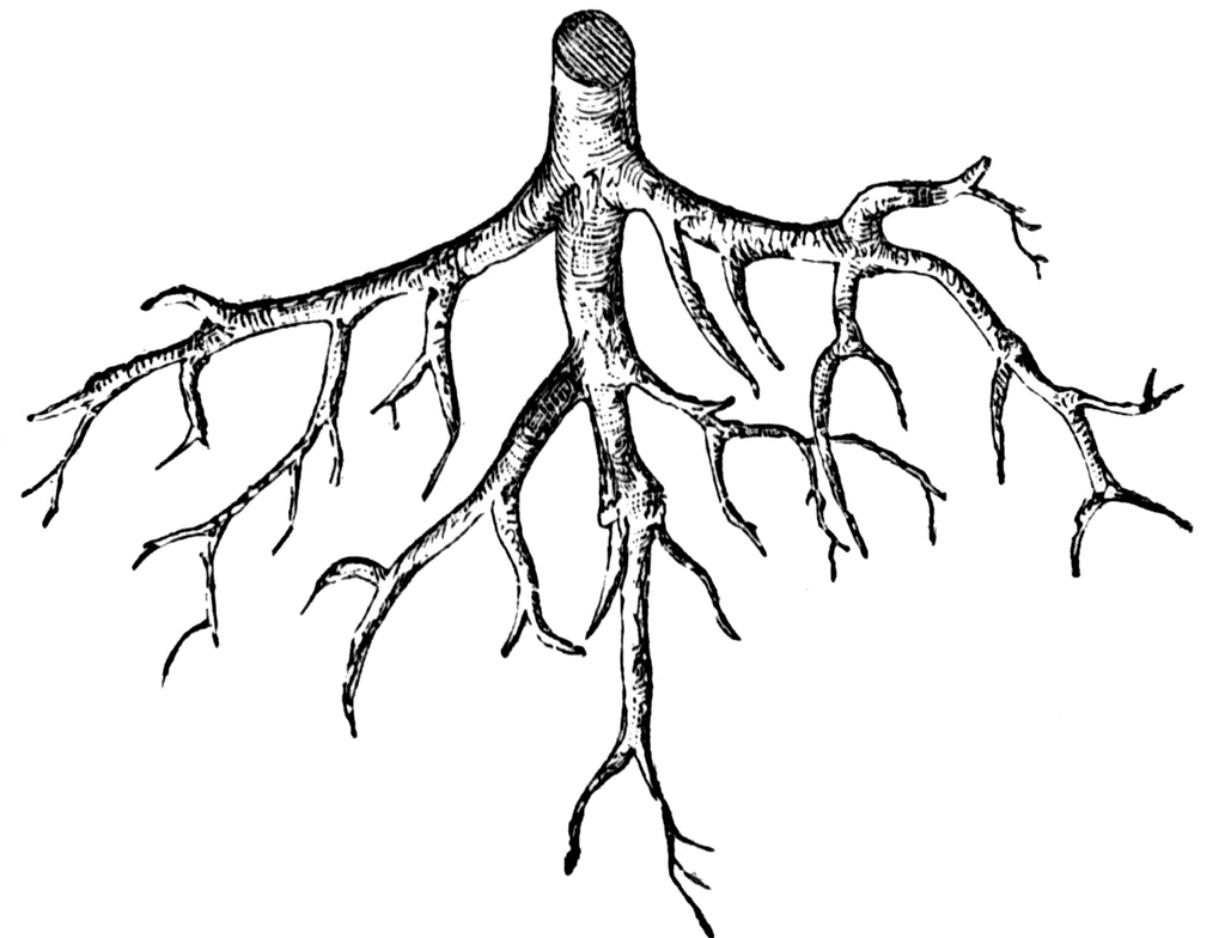
Can we learn from Linux?





Study Linux main repo

Not a single (open) branch in Linus repo



Challenge 1

- Personal repos are beyond reach
- Local commits might never be observable

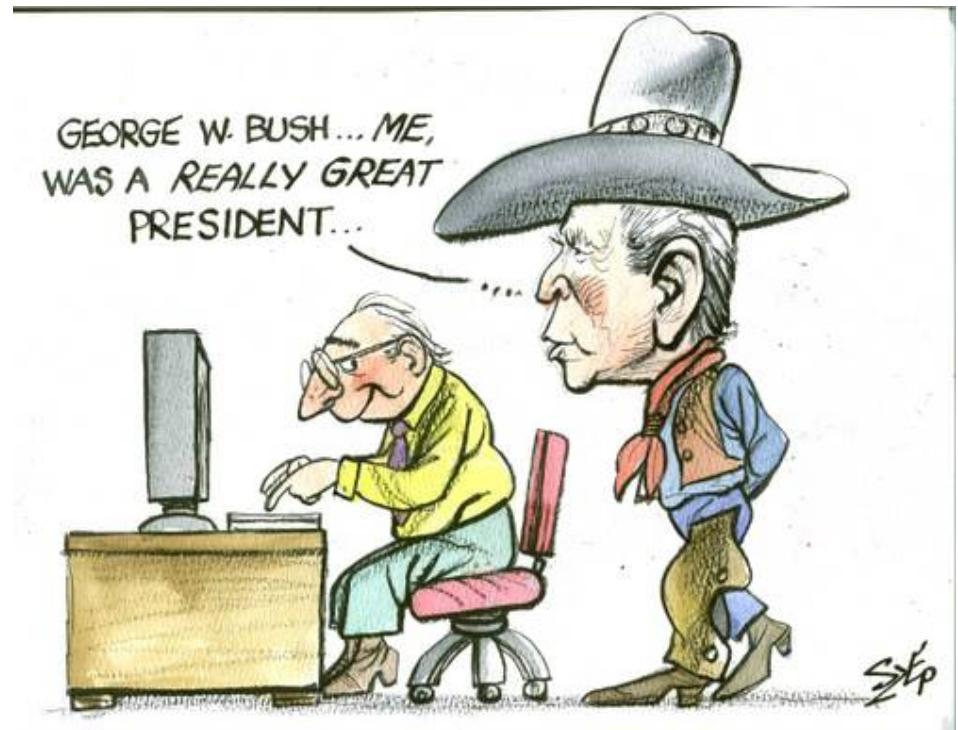


Challenge 2: History

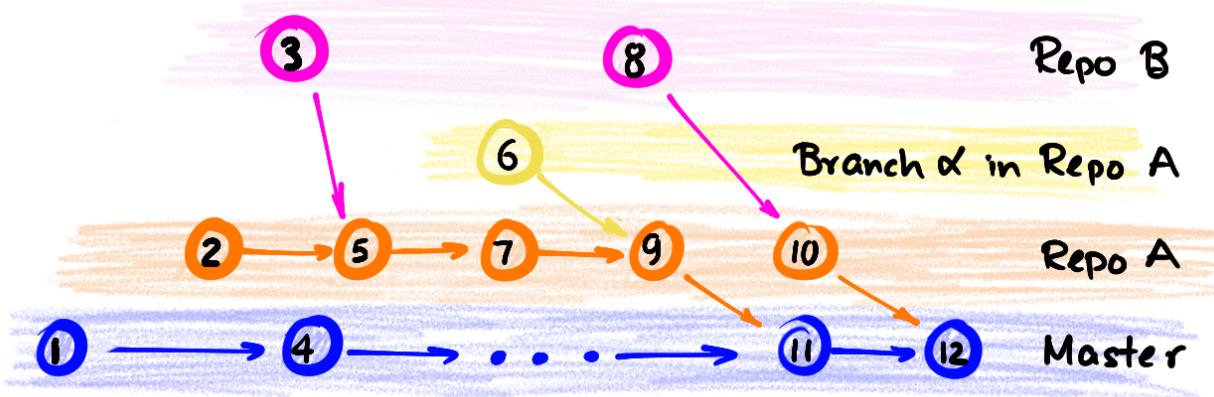
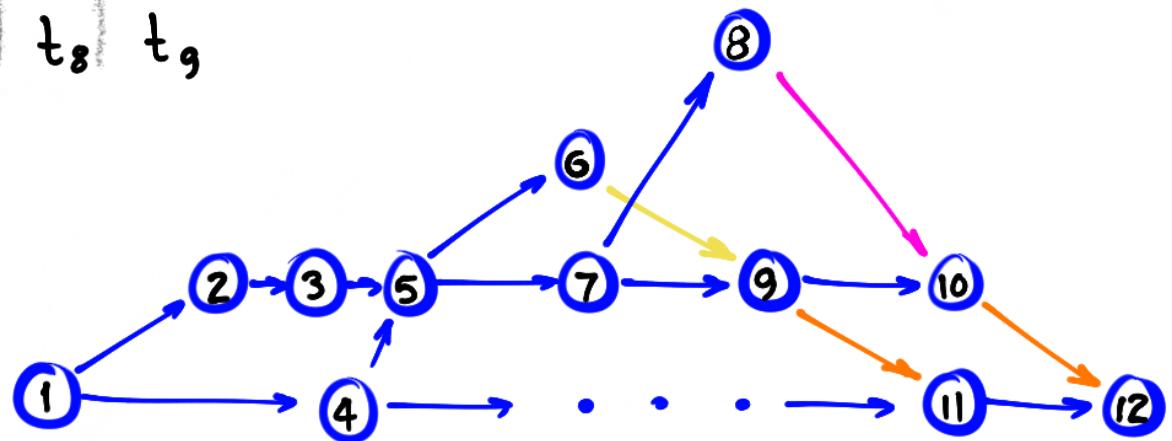
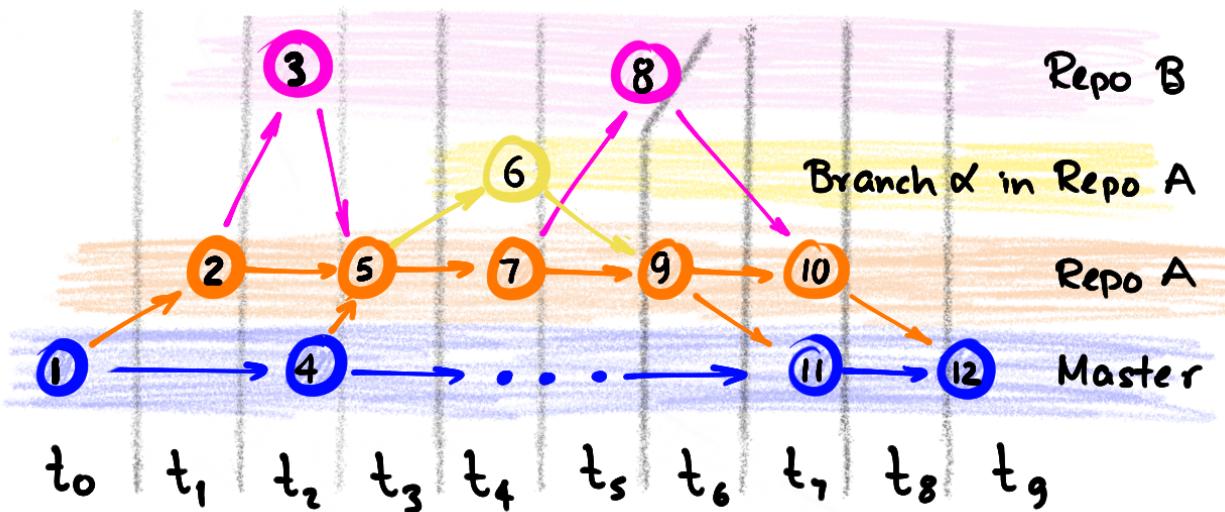


“History is written by the victors”

- Niccolò Machiavelli



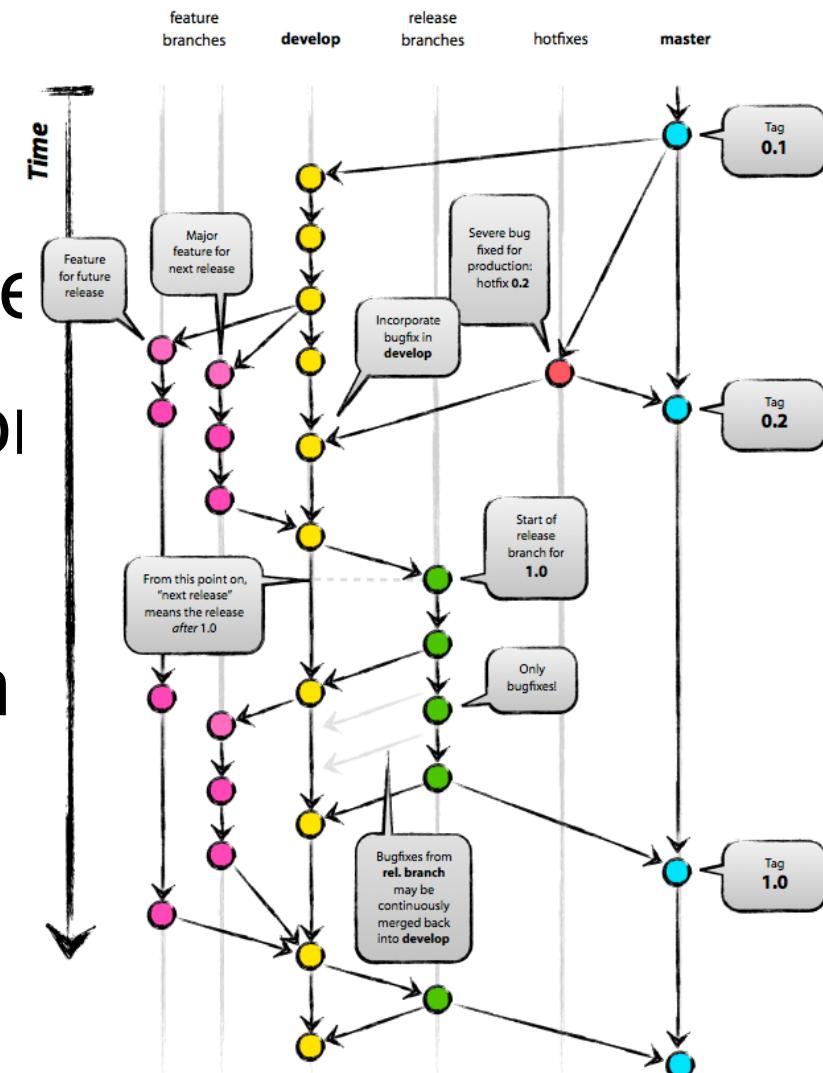
Save history before it is lost!



Super-repository

- Collection of repositories cloned (recursively) from the same repo

- At least one per developer
 - In their personal computer
- At least one public repository
 - The **blessed**
- In git, no way to trace them

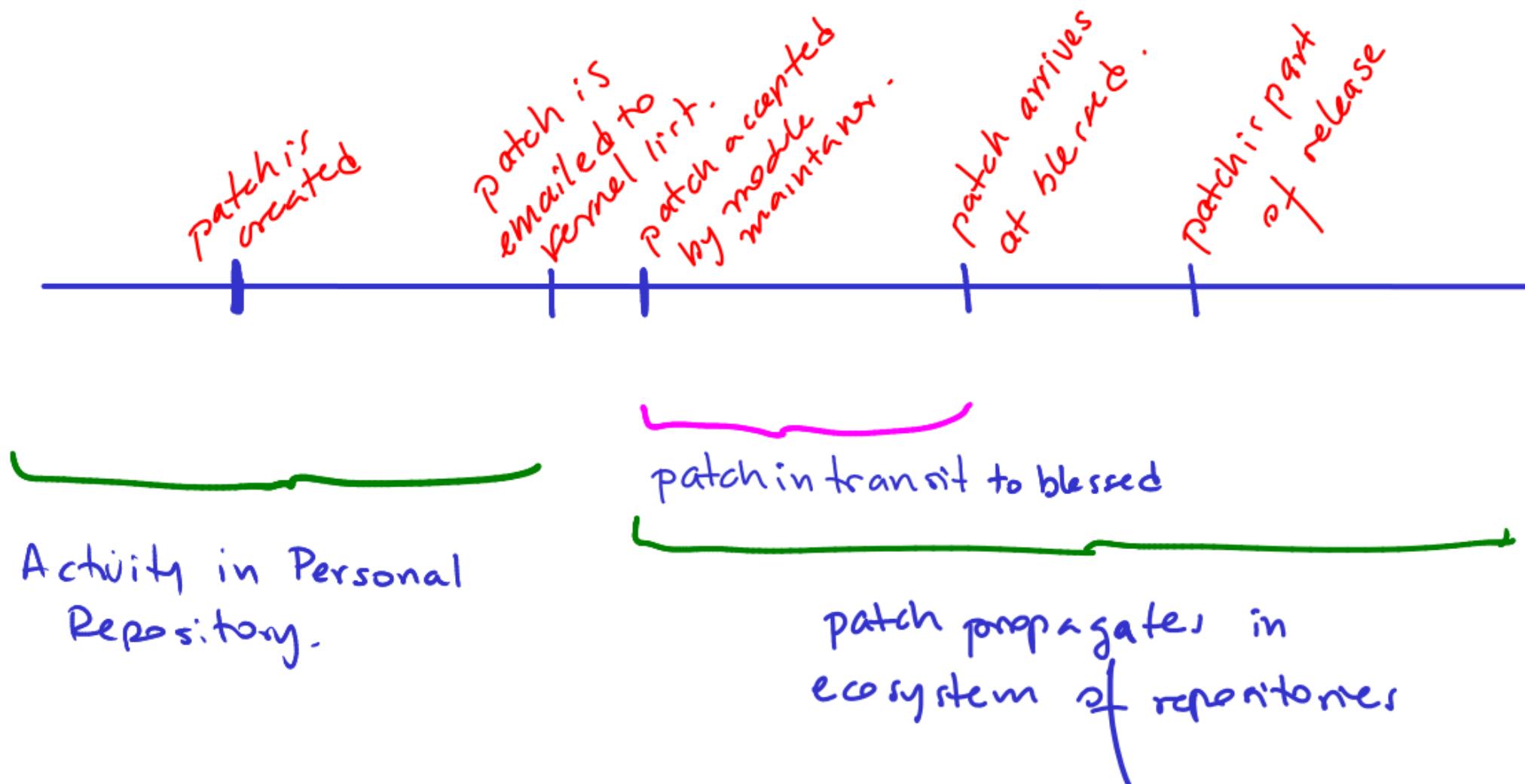


Moving commits across the superRepo

Method	
Push	Done at source, needs write access to source
Pull	Done at destination, needs read access to source
Email	Source creates patch, recipient applies it



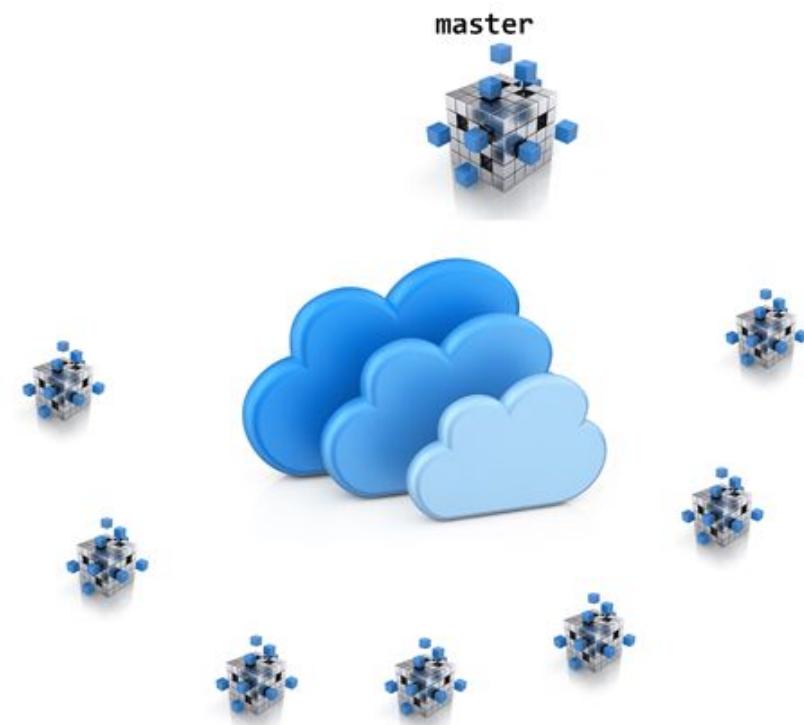
Life of a Patch in Linux



How can we observe them?

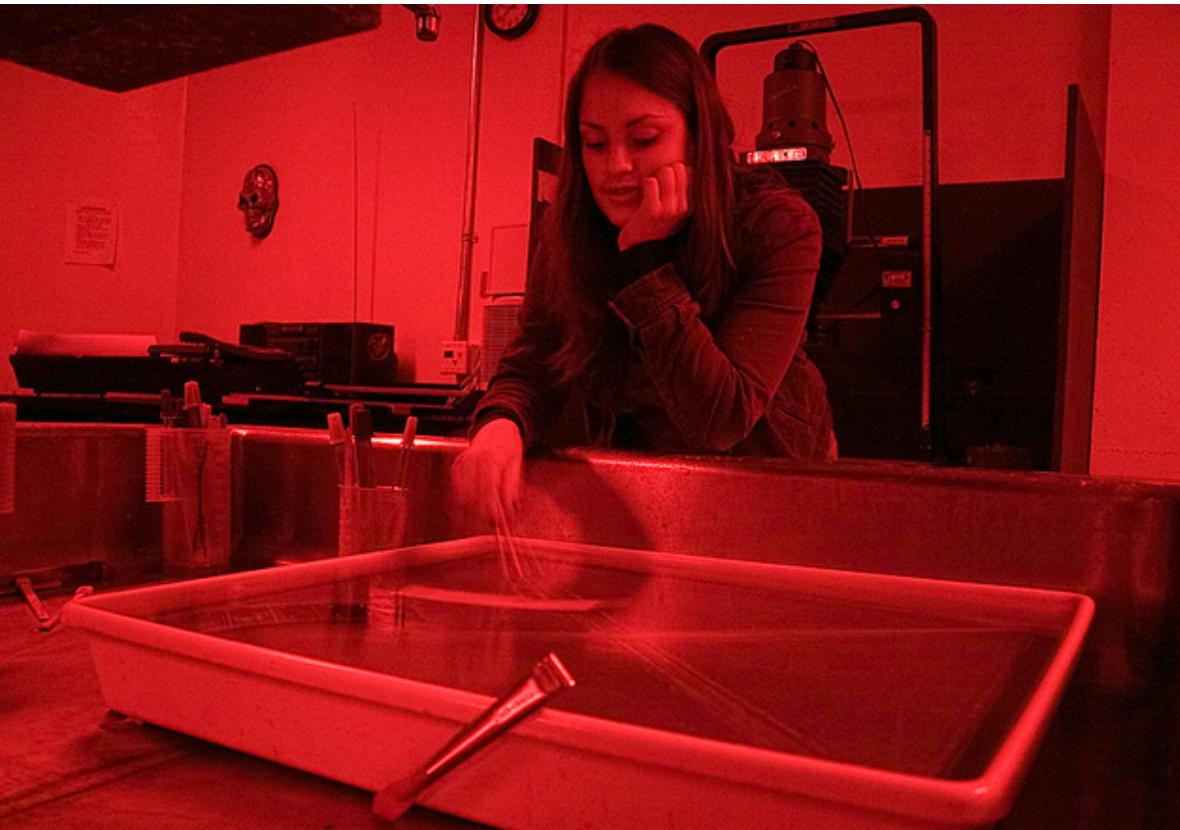


We have to find them!





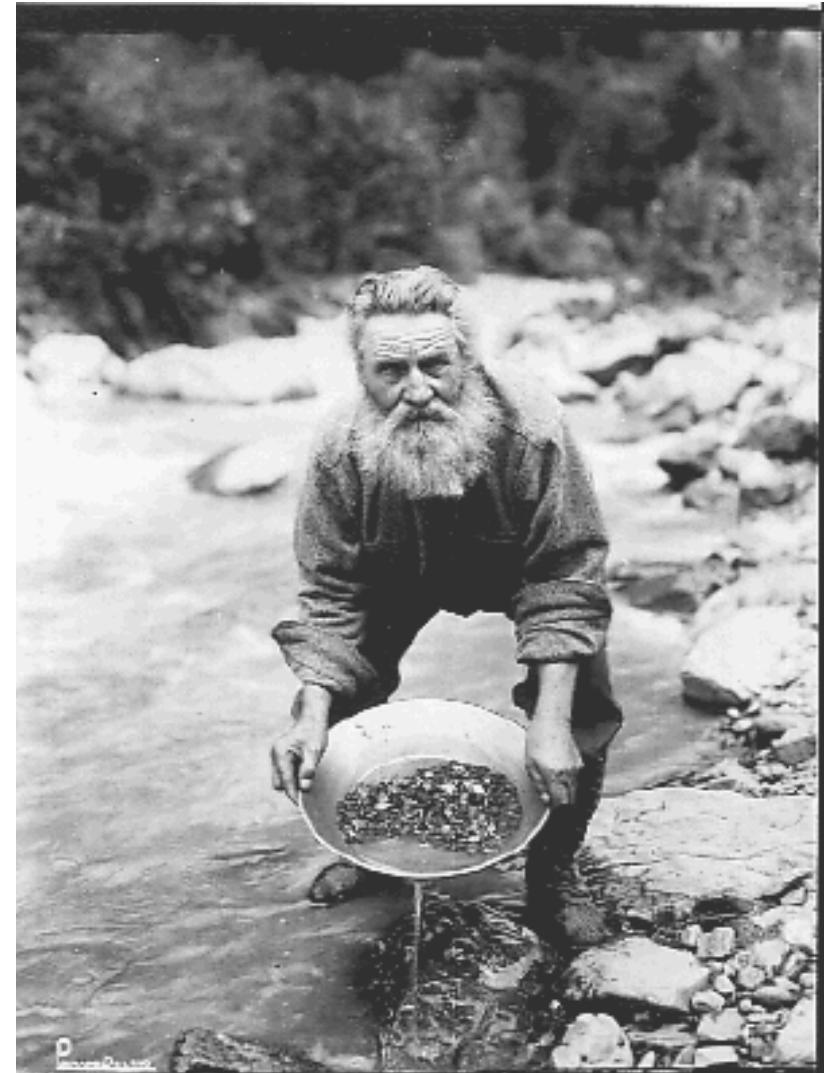
Snapshot Mining



Continuous Mining



Continuous Mining



- Every interval
 - Mine as many repos as known
 - What is new?
 - What has been deleted?

Continuous Mining

- Challenges:

- 1) Finding the repositories
- 2) Logging their commits

- 1) Proactive vs Reactive implementation

- Logging vs discovery

ContinuousMining of Linux

- Linux has no centralized logging
 - Nobody really knows what the superRepo is
 - Commits flow without any event broadcasting mechanism
- Where do we find the activity?
 - Repos
 - Commits
 - Mailing lists



Repos and Committers

- Most repos will have a known set of persons committing to them
 - Simplest case: its owner is the only committer
 - Extreme case: repo is used as centralized version control system: everybody commits to it

Semiautomatic Process

- **Every 3 hrs**, ask every repo:

- What new commits do you have?
 - What commits did you delete?
 - Automatically resolve propagations
 - Commits might propagate before we scan

- Daily:

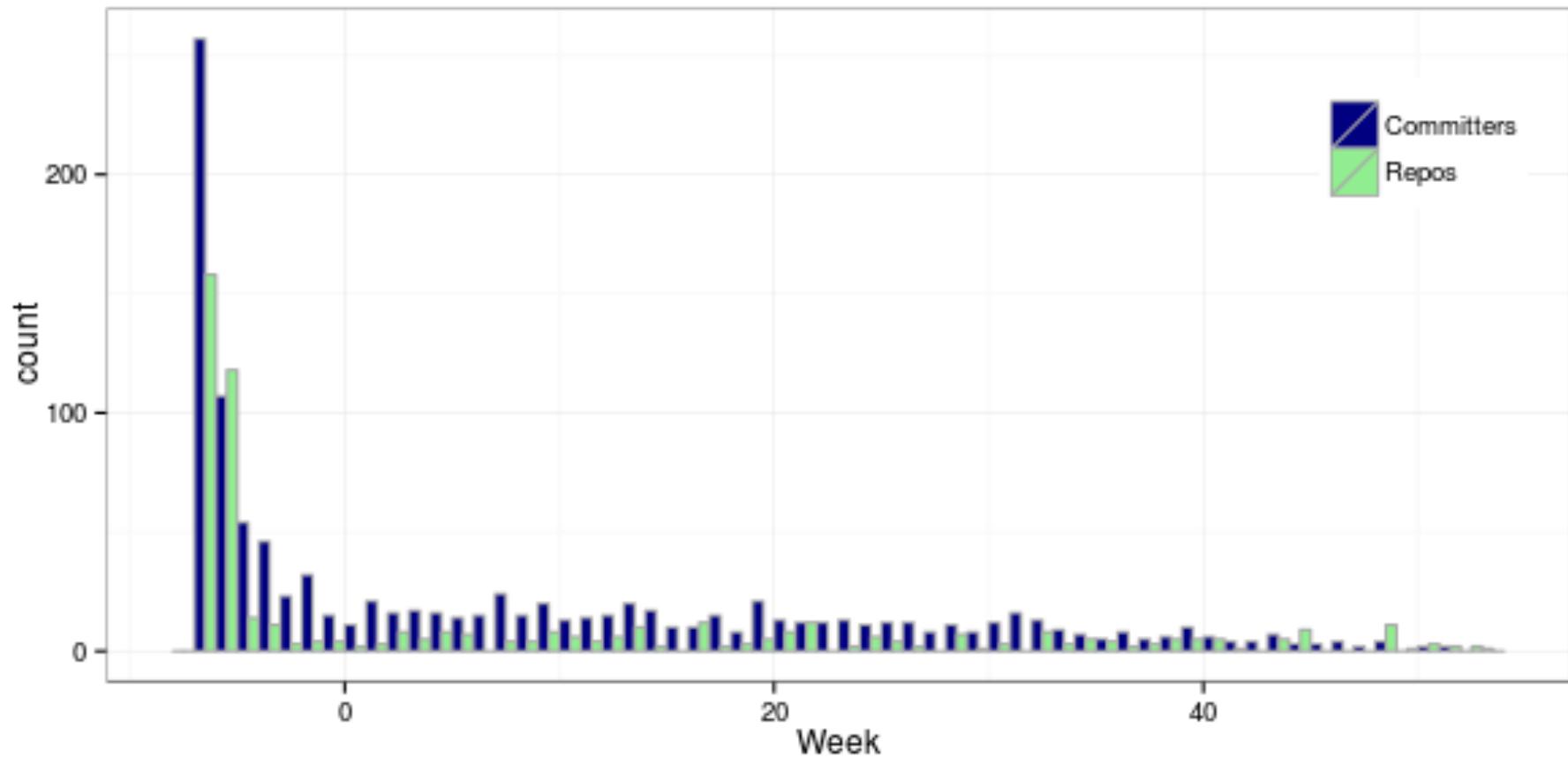
- Are commits in repo by unknown committers?
 - Answer:
 - **is there a new repo? or is committer new to repo?**

Implementation

- Run from Nov 2011 to April 2015
- Retrieved:
 - 3.8 M commits (compared to .55M in Linus repo)
 - 670 repositories
 - 160 million propagation records:
`<cid, added|deleted,repo,when>`



Discovery of new repos



Is one better than the other?

- RQ1

- Does continuousMining uncover a larger development ecosystem than snapMining?

- RQ2

- Does continuousMining expose any missing information, or bias in the recorded history of the project recovered using snapMining?



	Snapshot (Linus)	Continuous
No Repos	1	479
Commits	64k	533k
Non-merge Commits	59k	485k
Unique Non-merges	58k	135k
%unique non-merges	98.9%	27.9%
Non-merges that reached Blessed		43.1%
Different authors emails	3434	5646
Different authors	2883	4575
Different committers emails	283	1185
Different committers	245	1058

- RQ2

- Does continuousMining **expose any missing information**, or bias in the recorded history of the project recovered using snapMining?

Commit vs Patches

	Blessed	Super- repository
Non-merge commits	58,953	485,027
Patches	58,356	135,532
Ratio of patches to commits	98.9%	27.9%
Ratio of commits committed in 2012 that reached <i>blessed</i>		12.1%
Ratio of patches created in 2012 that reached <i>blessed</i>		43.1%

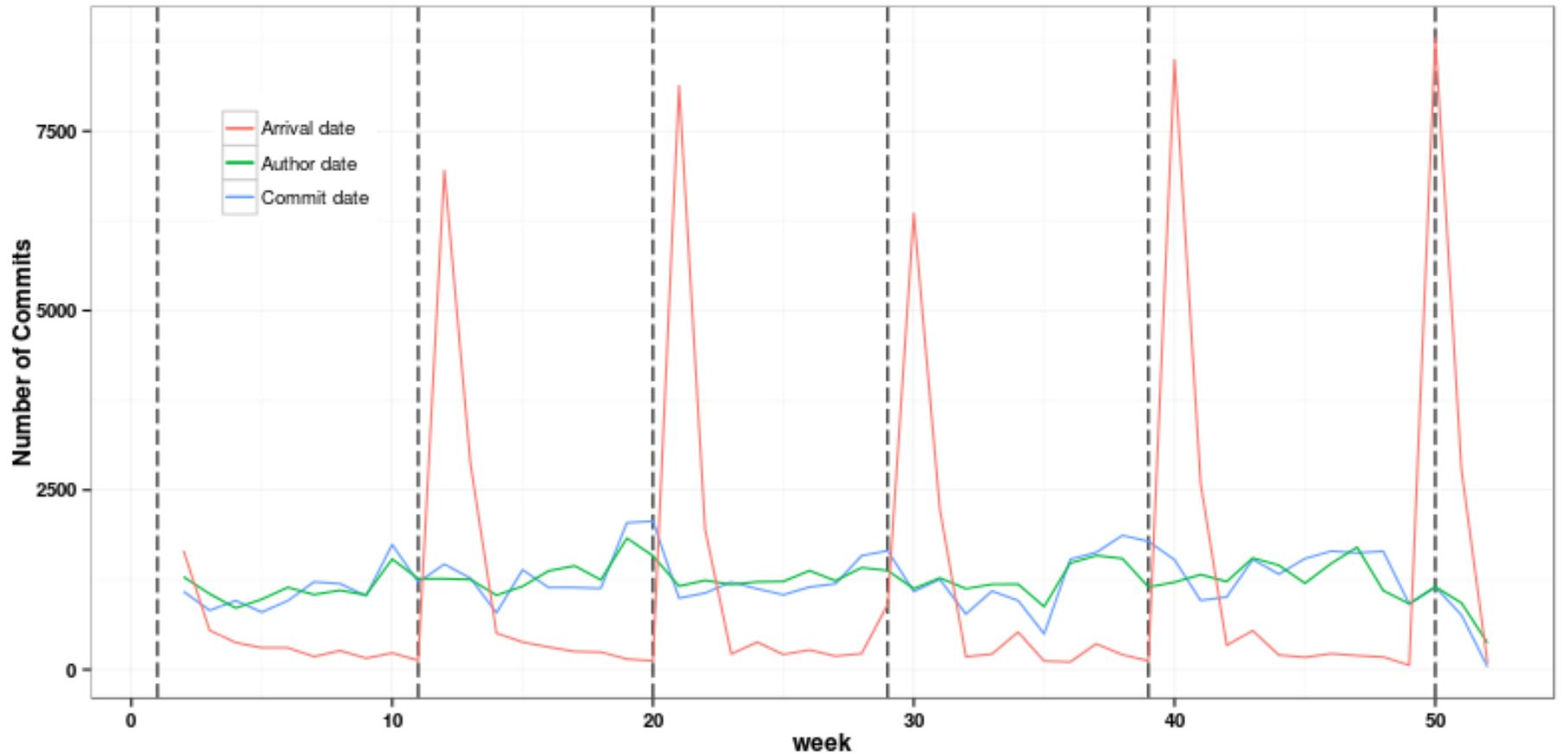
- Due to rebasing, commit **ids are insufficient** to tracks patches
- Large amount of work not reaching *blessed*

Committer data in *blessed* is biased

- 37.9% of patches arriving at ***blessed*** did not arrive in its original commit

	Count	Prop
Changed date of commit	22,326	37.9%
Changed committer	7,330	12.4%
Changed date of authorship	6,725	11.4%
Changed author	203	0.3%
Changed log	13,323	22.6%
Changed log summary (subset of change log)	2,634	04.5%

Arrival of Commits at Blessed



patch is created
Patch is emailed to kernel list.
Patch accepted by maintainer.
Patch arrives at blessed

Arrival of Commits at Blessed...

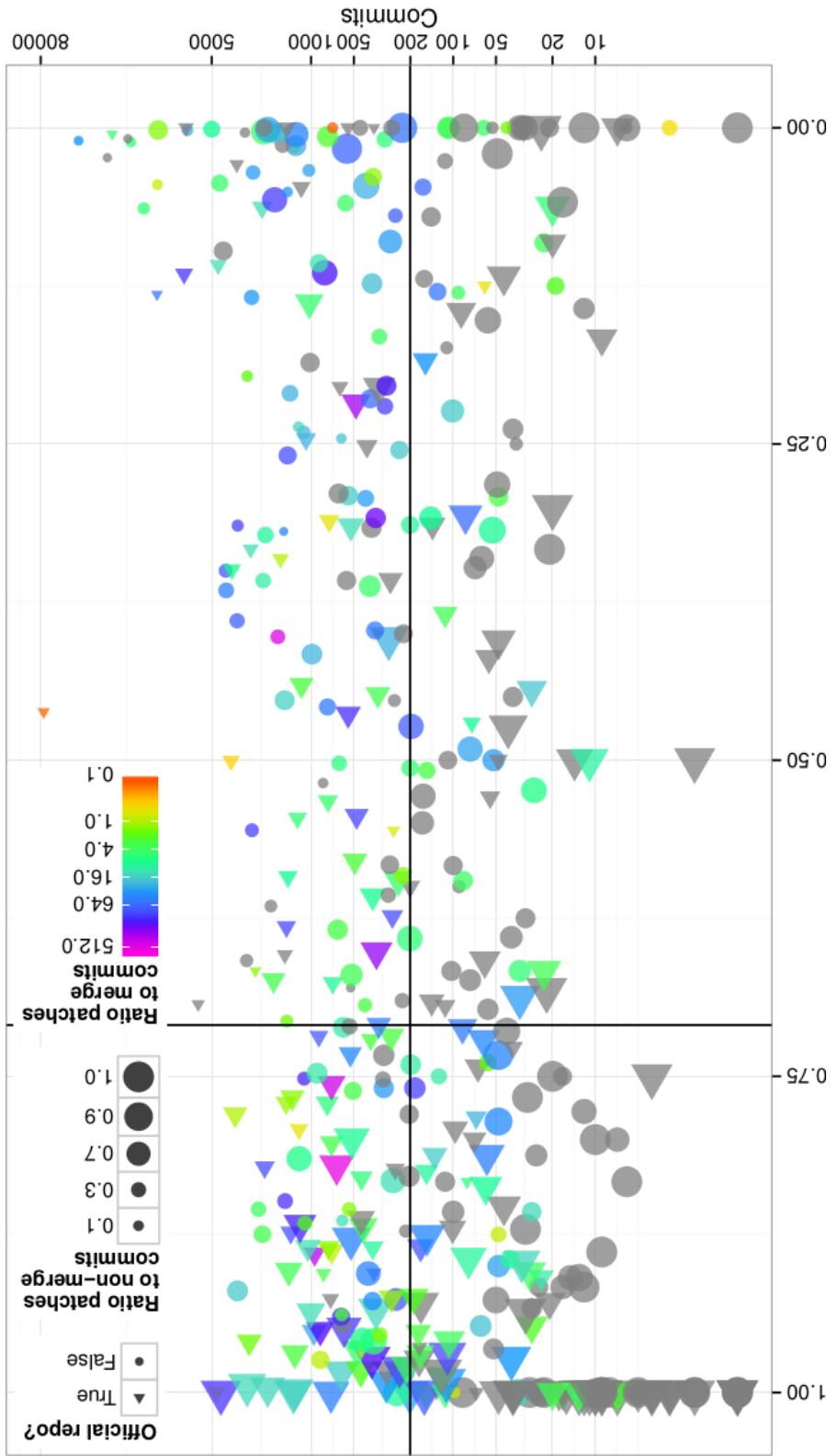
- We can classify patches as a new feature or bug-fix

Release	Merge Window	Ratio Identified Bugs	Bug-Fixing Window	Ratio Identified Bugs	Total
3.2	9,459	15.8%	1,955	42.9%	11,415
3.3	9,858	17.5%	1,969	41.2%	11,828
3.4	10,068	13.3%	1,702	42.1%	11,771
3.5	9,160	15.1%	1,968	41.9%	11,129
3.6	11,083	15.4%	1,920	42.4%	13,004
3.7	11,737	15.4% (incomplete)			

Empirical study

- RQ. What are the characteristics of the repos in the Linux Super-repository

	Min	1st Q	Median	Mean	3rd Q	Max	<i>blessed</i>	Total	<i>Snap.</i>
Total Commits	1	34	198	1,162.3	717.5	76,279	4,633	533,513	64,029
Contributed to <i>blessed</i>	0	2	23	139.5	117.5	4633	4633	64,029	
Ratio contributed to total	0	1%	29%	38%	67%	100%	100%		
Merges	0	0	2	105.6	19.5	23,078	2,103	48,486	5,076
Non-merges	0	33	169	1056.7	617.5	53,201	2,530	485,027	58,953
Patches	0	18	84	269.1	286.5	5,911	1,116	125,535	58,356
Committees	1	1	1	3.2	1	339	1	1,058	245
Authors	1	3.5	12	44.51	38.5	1,061	454	4,575	2,883
Ratio Authors/ Committees	1	4	10	25.5	28.3	454	454	4.32	11.76



The Repos

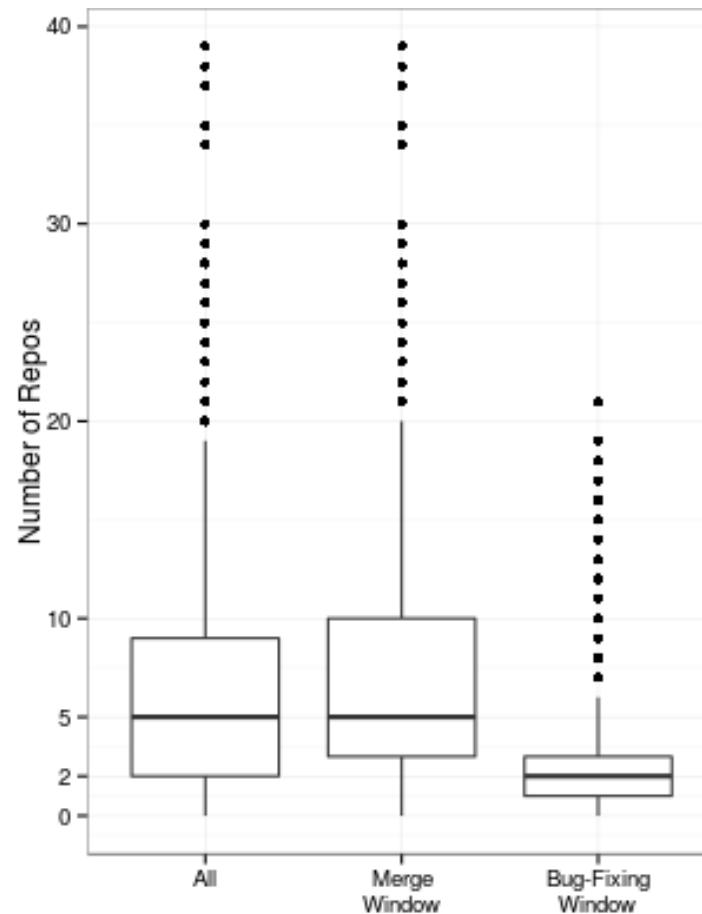
- X: activity (in commits)
- Y: ratio of commits accepted by Linus to total commits
- Shape:
 - Triangles: official repos
 - Circles: non-official repos
- Size:
 - Smaller: consume commits
 - Larger: produce commits
- Color: merge/commit ratio
 - Grey: never merge
 - “Cooler”: high ratio
 - “Warmer”: lower ratio

Consumers are very active

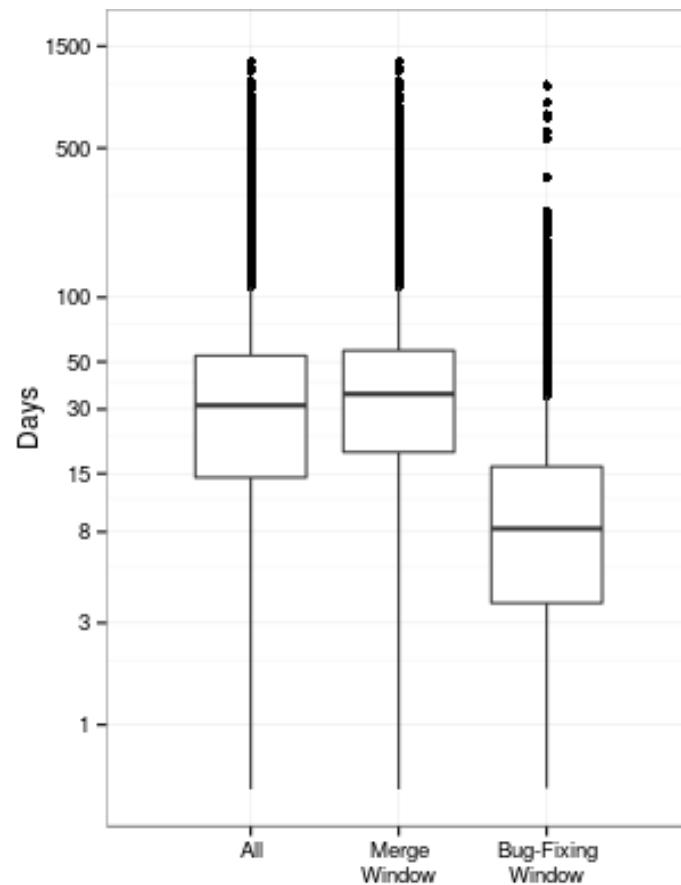
Address	New Com- mits	Com- mit- ters	New Patch.	Contr. Patch.	Prop. Patches Contr.
linaro.org/.../andygreen/kernel-tilt	43,290	21	2743	28	1.0%
ubuntu.com/.../ubuntu-precise	27,141	6	1184	28	2.3%
kernel.org/.../rt/linux-stable-rt	25,129	1	407	2	0.4%
ubuntu.com/.../ubuntu-quantal	19,459	6	809	7	0.8%
linaro.org/.../linux-tb-ux500	18,508	21	1765	20	1.1%
linaro.org/..../linux-linaro-tracking	12,058	28	1291	58	4.4%
android.googlesource.com/.../msm	11,922	339	5911	10	0.1%
kernel.org/.../stable/linux-stable	7,591	3	1116	0	0.0%
github.com/vstehle/linux	7,421	14	1079	2	0.1%
denx.de/linux-denx	5,005	69	1915	2	0.1%
opensuse.org/kernel	4,399	41	1716	75	4.3%

Propagation

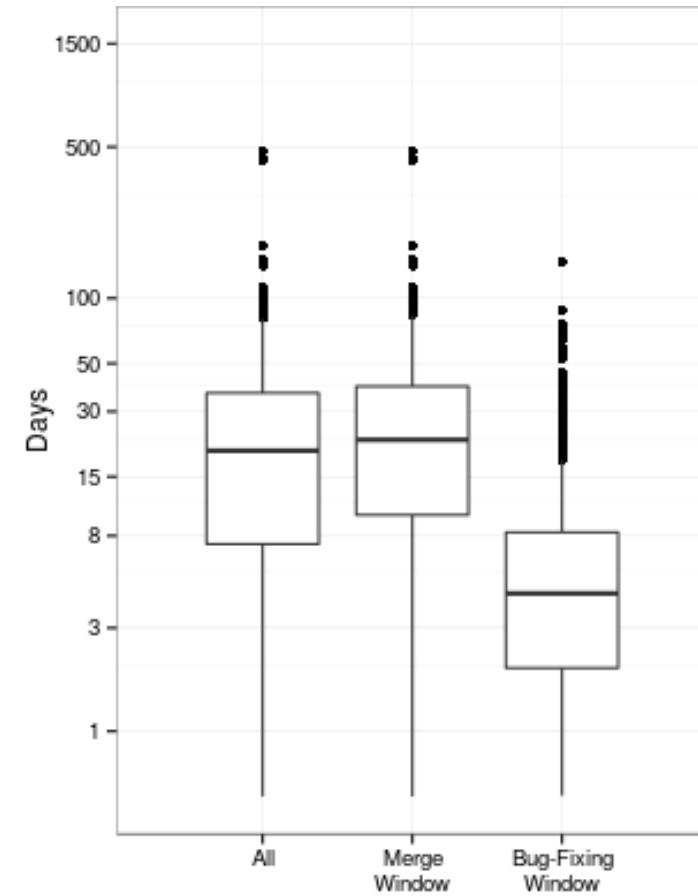
- RQ: How do repositories interact, and how do commits propagate across repositories?



The Latency

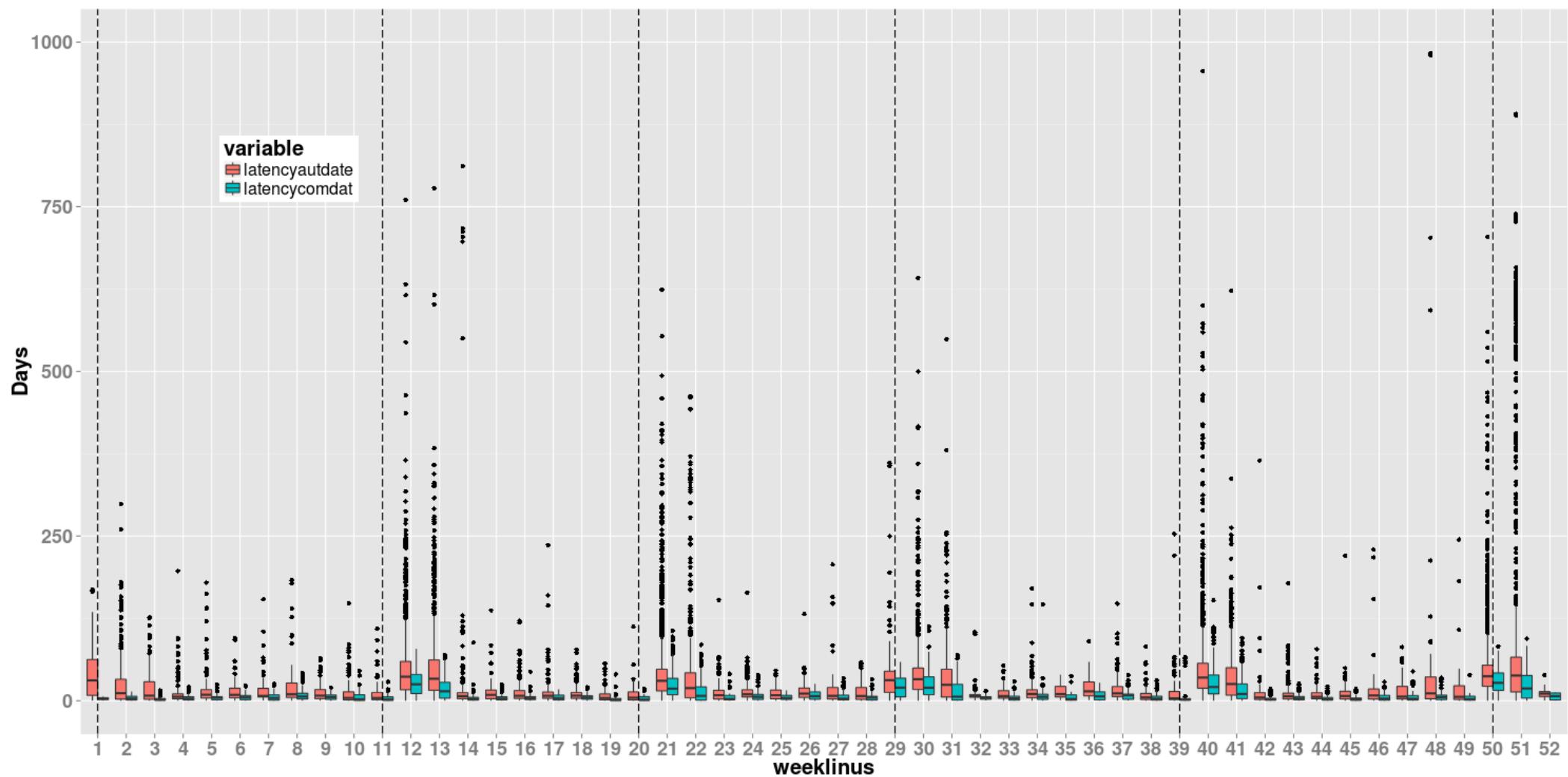


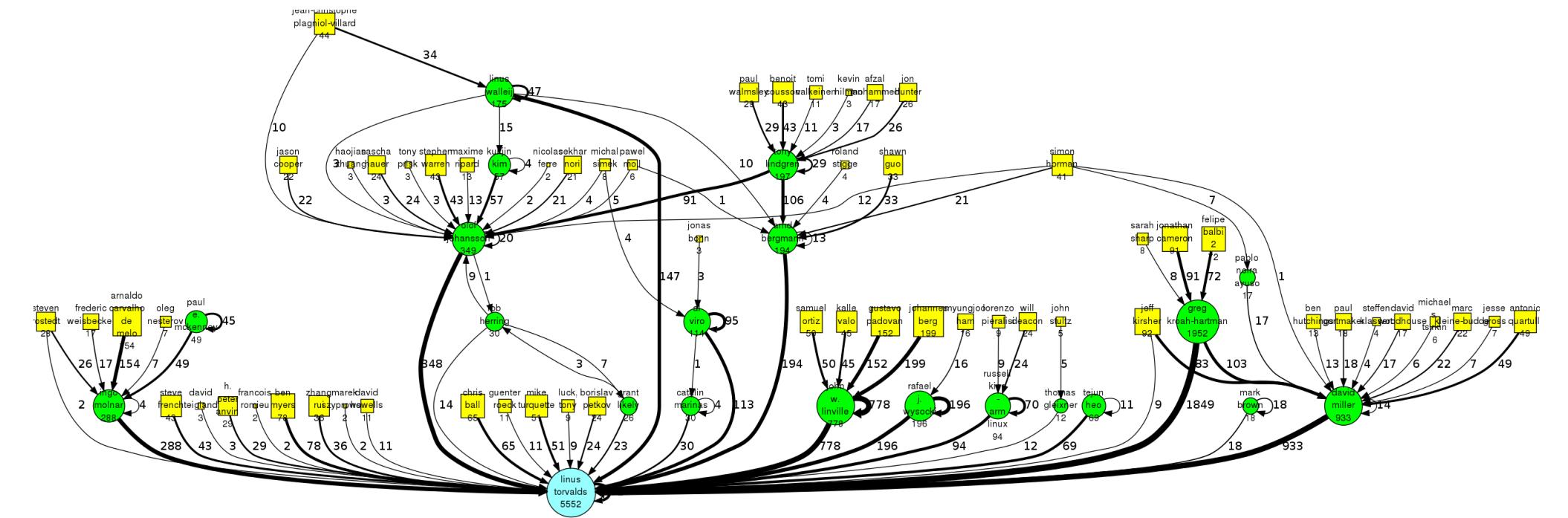
Time of Authorship



Time of Commit

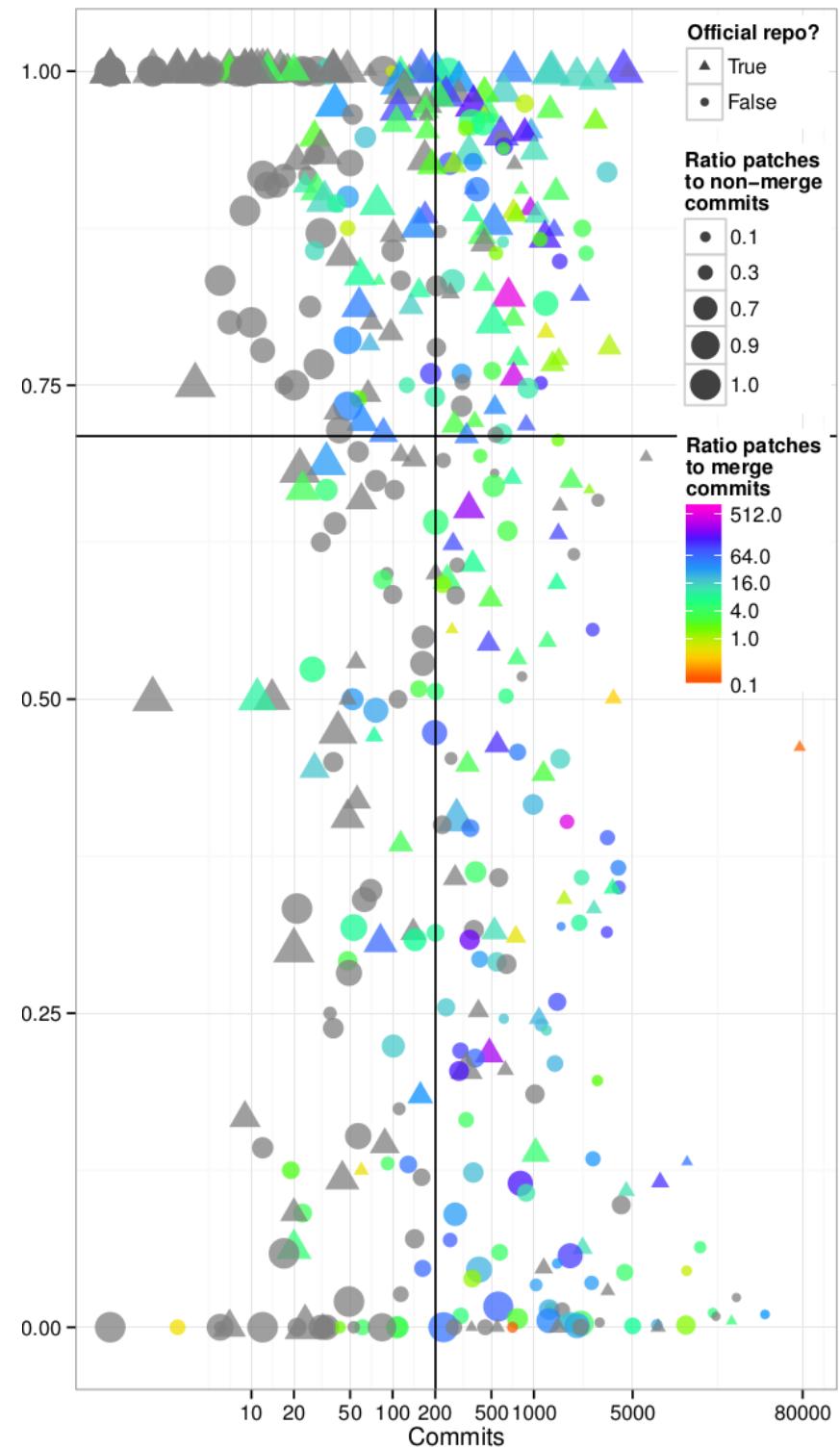
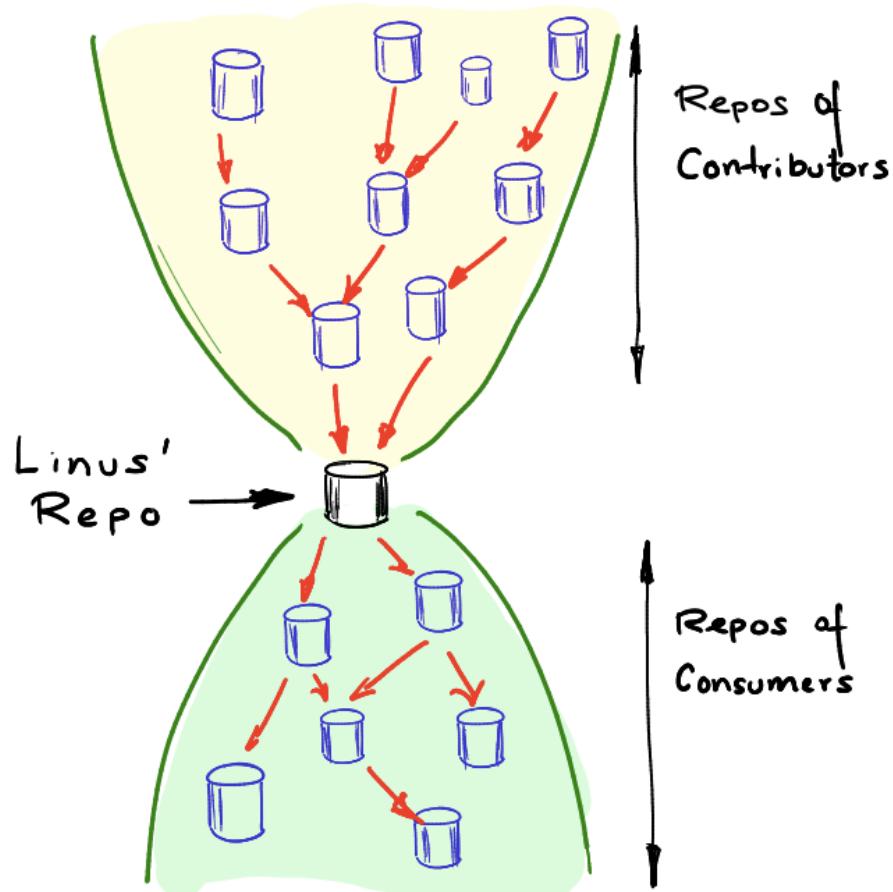
Latency





The Linux Super-repository

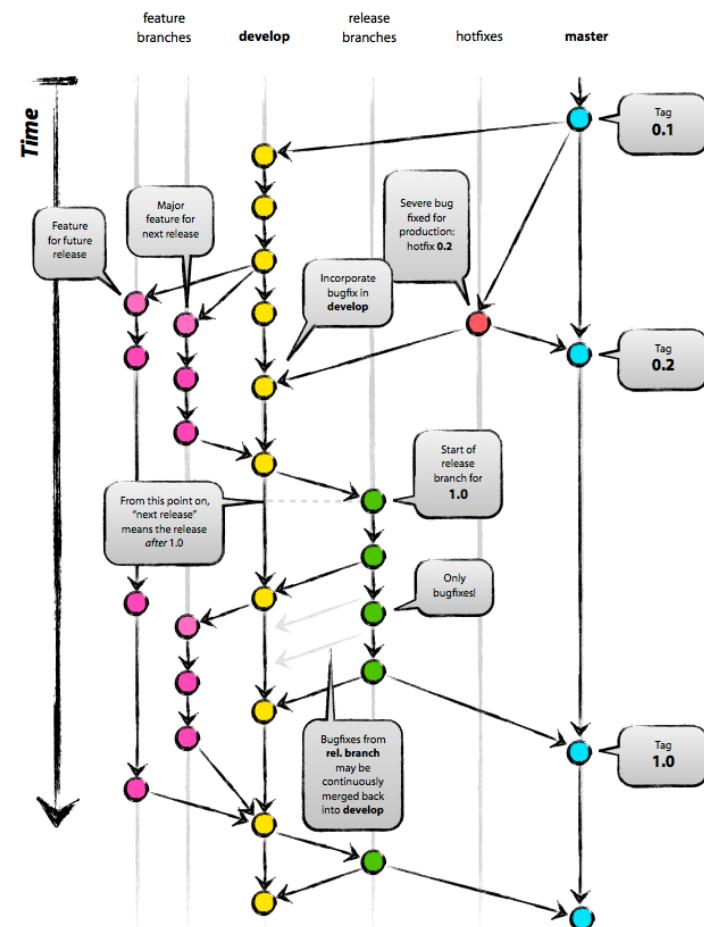
- Producers
- Consumers



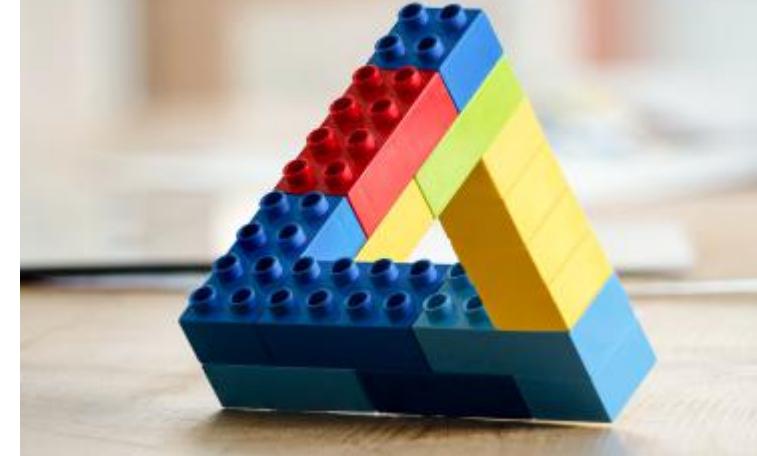
Final observations

Super-repository

- Collection of **repositories cloned (recursively)** from the same repo



The paradox



**Distributed Version control allows
anybody to participate...**

**but it centralizes control on who
contributes**

Linux is not only software...

it is also software engineering

Git promotes Linux practices:

Git promotes Linux practices:

- code reviews**

Git promotes Linux practices:

- code reviews**
- pull requests**

Git promotes Linux practices:

- code reviews**
- pull requests**
 - isolation**

Git promotes Linux practices:

- code reviews**
- pull requests**
- isolation**

and github popularizes them

GitHub for non-technical people

GitHub



GitHub for the rest of us

contributions welcome

ROBERT MCMILLAN BUSINESS 09.02.13 6:30 AM

FROM COLLABORATIVE CODING TO WEDDING INVITATIONS: GITHUB IS GOING MAINSTREAM



Credit: Flickr/Nick Quaranto

Git made it possible for programmers to coordinate distributed work across teams -- now GitHub makes it possible for everyone else