

Principles of Software Construction: Objects, Design, and Concurrency

Software development at scale

Bonus slides: Unseen GoF design patterns

(The end)

Michael Hilton

Bogdan Vasilescu

Administrivia

- Final exam Monday May 6th 5:30-8:30 GHC 4401
- Review session Saturday May 4th 1pm NSH 3305

Intro to Java

Git, CI

UML

Static Analysis

GUIs

More Git

Streams

**Software
Engineering
in Practice**

Design



Part 1:

Design at a Class Level

**Design for Change:
Information Hiding,
Contracts, Unit Testing,
Design Patterns**

**Design for Reuse:
Inheritance, Delegation,
Immutability, LSP,
Design Patterns**

Part 2:

Designing (Sub)systems

Understanding the Problem

**Responsibility Assignment,
Design Patterns,
GUI vs Core,
Design Case Studies**

**Design for Reuse at Scale:
Frameworks and APIs**

Part 3:

**Designing Concurrent
Systems**

**Concurrency Primitives,
Synchronization**

**Designing Abstractions for
Concurrency**

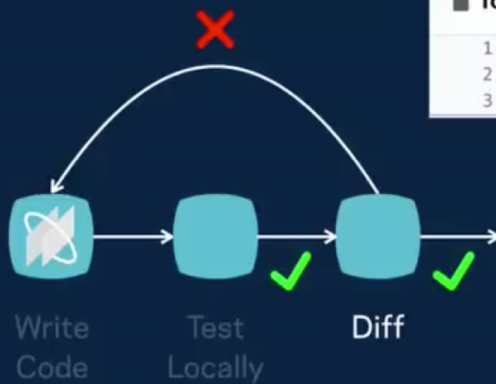
SOFTWARE DEVELOPMENT AT SCALE

Releasing at scale in industry

- Facebook: <https://atscaleconference.com/videos/rapid-release-at-massive-scale/>
- Google: <https://www.slideshare.net/JohnMicco1/2016-0425-continuous-integration-at-google-scale>
 - <https://testing.googleblog.com/2011/06/testing-at-speed-and-scale-of-google.html>
- Why Google Stores Billions of Lines of Code in a Single Repository: <https://www.youtube.com/watch?v=W71BTkUbdqE>
- F8 2015 - Big Code: Developer Infrastructure at Facebook's Scale: <https://www.youtube.com/watch?v=X0VH78ye4yY>

Pre-2017 release management model at Facebook

Diff lifecycle: First, local testing



Tools/xctool/xctool/xctool/Version.m View Options

1 #import "Version.h"

2

3 NSString * const XCToolVersionString = @"0.2.1";

1 #import "Version.h"

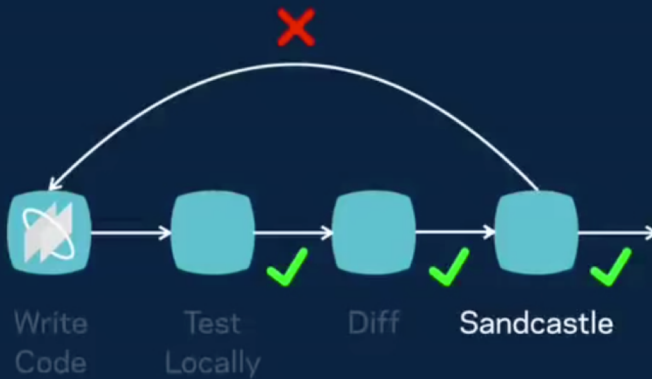
2

3 NSString * const XCToolVersionString = @"0.2.2";

```
PASS ExampleTest (0.050s)
.
OK (1 test, 4 assertions)
OK
(1 tests, 4 assertions, 0 incomplete, 0 failures)
```

Test and lint locally

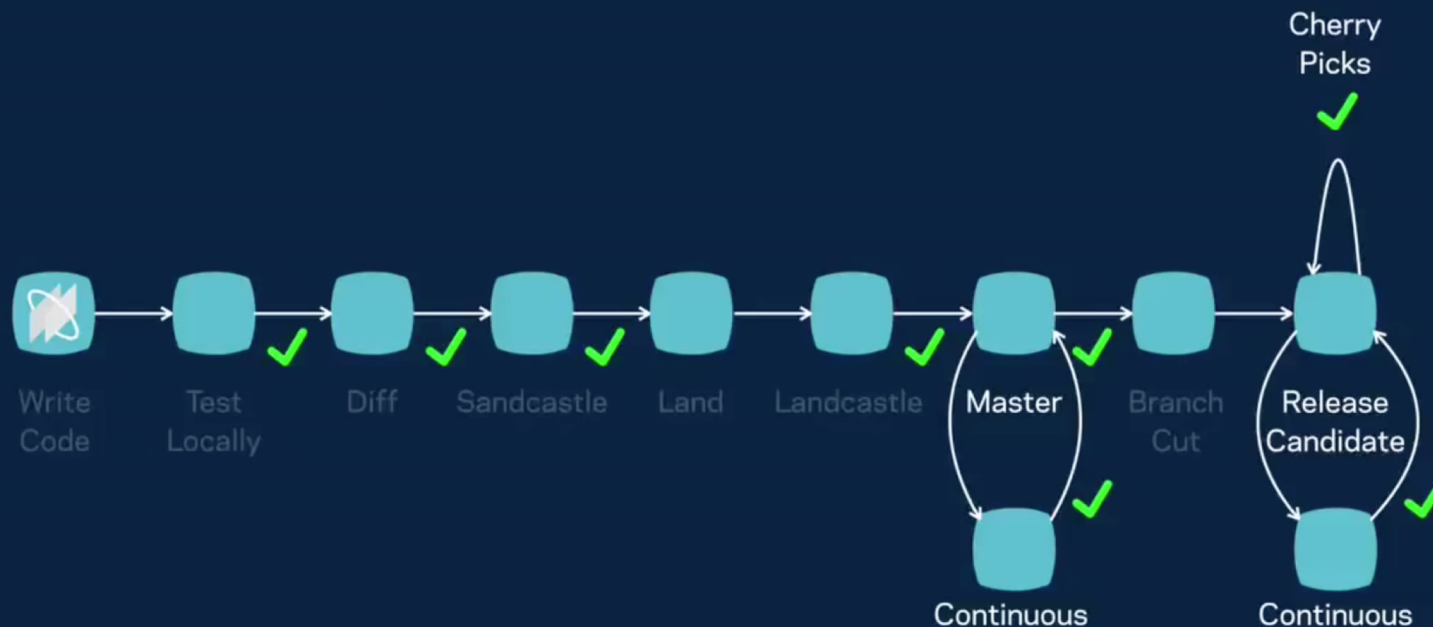
Diff lifecycle: Next, CI testing (data center)



	Facebook	Messenger	Groups	...
arm	✓	✓	✓	✓
x86	✓	✓	✓	✓
...	✓	✓	✓	✓

App and Build
Configuration Matrix

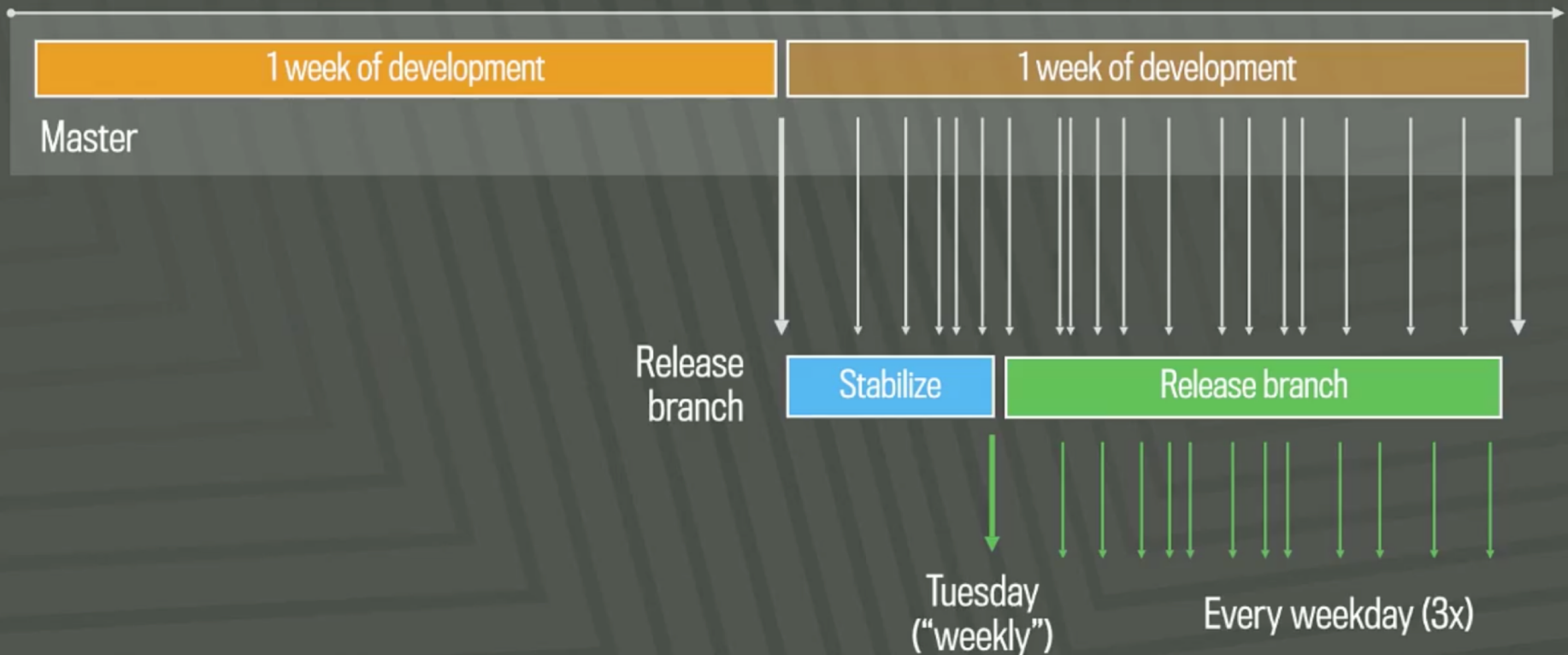
Diff lifecycle: Then, diff ends up on master



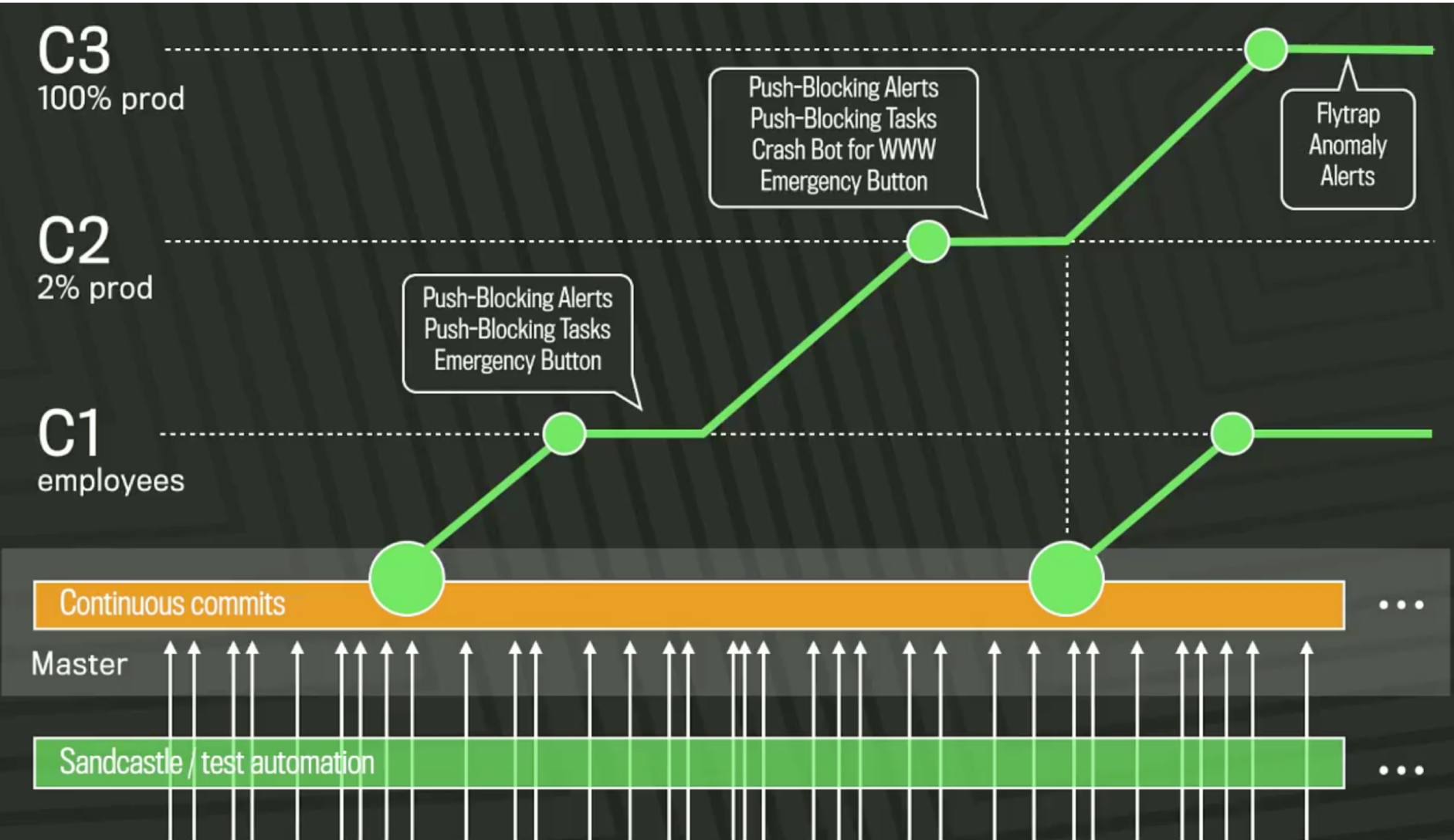
Dogfooding

Release every two weeks

www.facebook.com



Quasi-continuous push from master (1,000+ devs, 1,000 diffs/day); 10 pushes/day



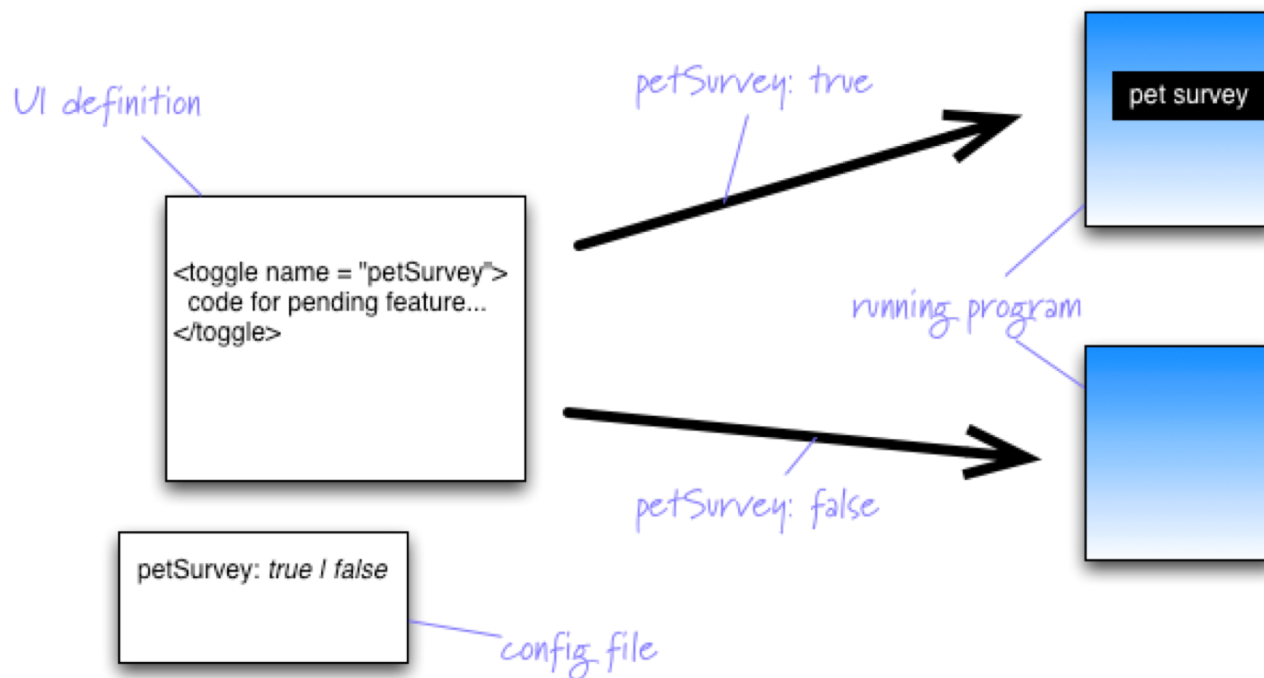
Aside: Key idea – fast to deploy, slow to release

Dark launches at Instagram

- **Early:** Integrate as soon as possible. Find bugs early. Code can run in production about 6 months before being publicly announced (“dark launch”).
- **Often:** Reduce friction. Try things out. See what works. Push small changes just to gather metrics, feasibility testing. Large changes just slow down the team. Do dark launches, to see what performance is in production, can scale up and down.
"Shadow infrastructure" is too expensive, just do in production.
- **Incremental:** Deploy in increments. Contain risk. Pinpoint issues.

Aside: Feature Flags

Typical way to implement a dark launch.



<http://swreflections.blogspot.com/2014/08/feature-toggles-are-one-of-worst-kinds.html>

<http://martinfowler.com/bliki/FeatureToggle.html>

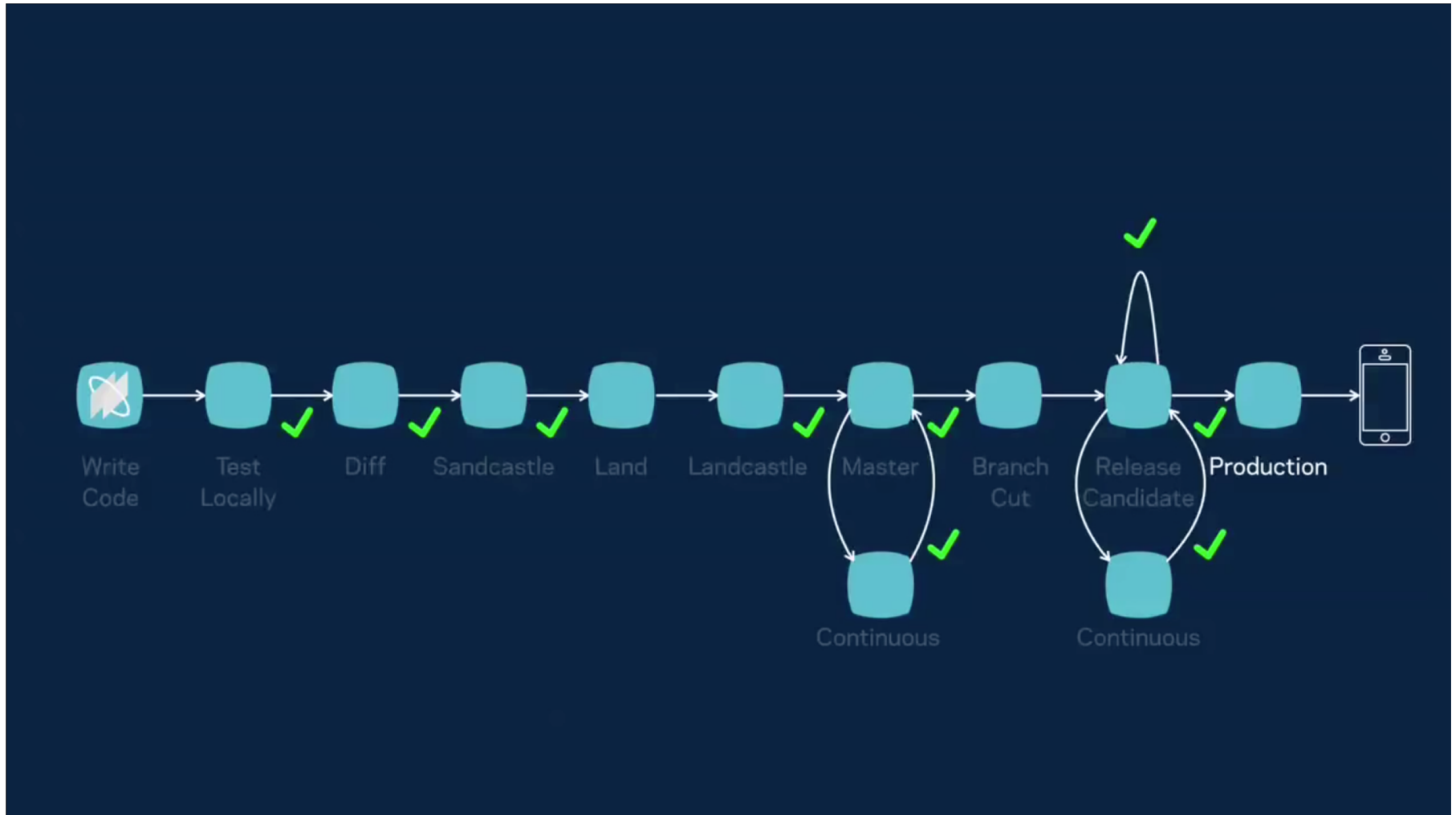
Issues with feature flags

Feature flags are “technical debt”

Example: financial services company with nearly \$400 million in assets went bankrupt in 45 minutes.

<http://dougseven.com/2014/04/17/knightmare-a-devops-cautionary-tale/>

Diff lifecycle: Finally, in production

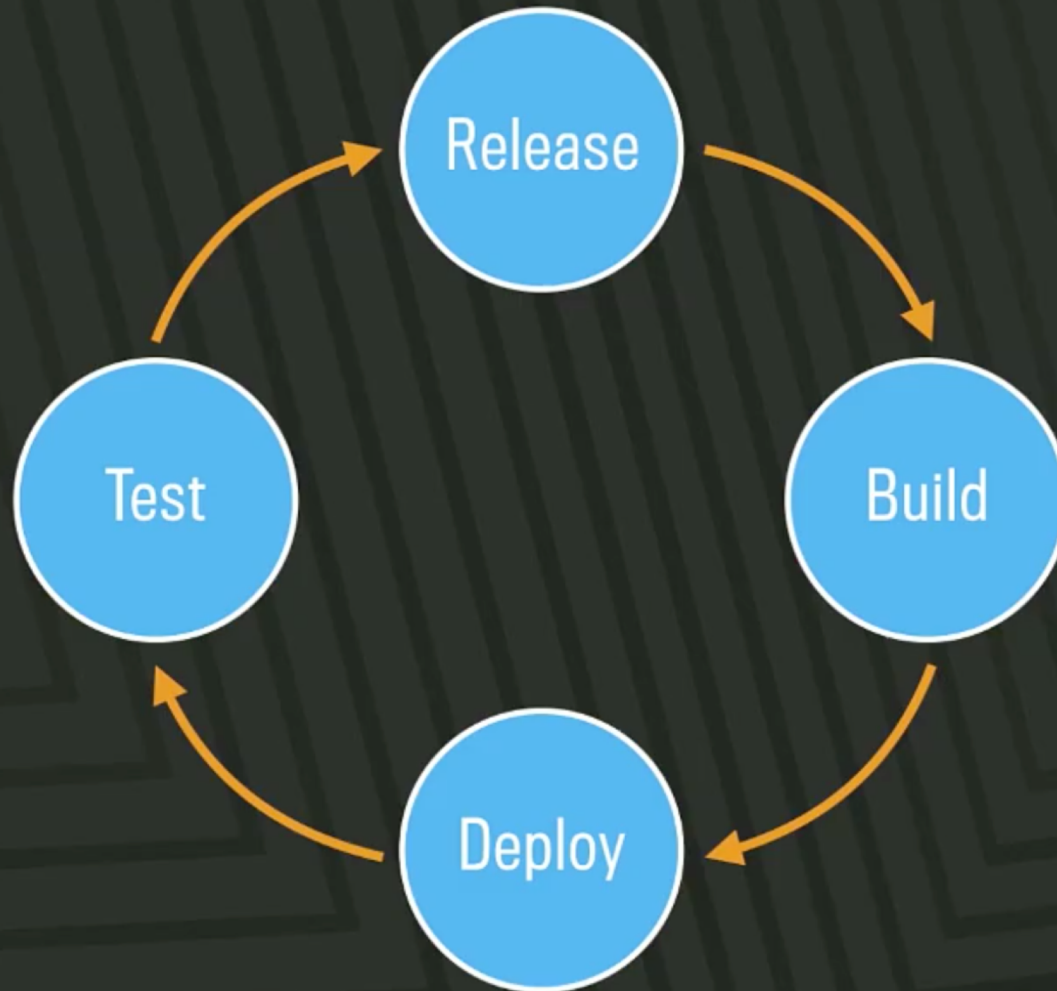


What's in a weekly branch cut? (The limits of branches)

Weekly web branch



Quasi-continuous web release



Google: similar story. HUGE code base

Google repository statistics

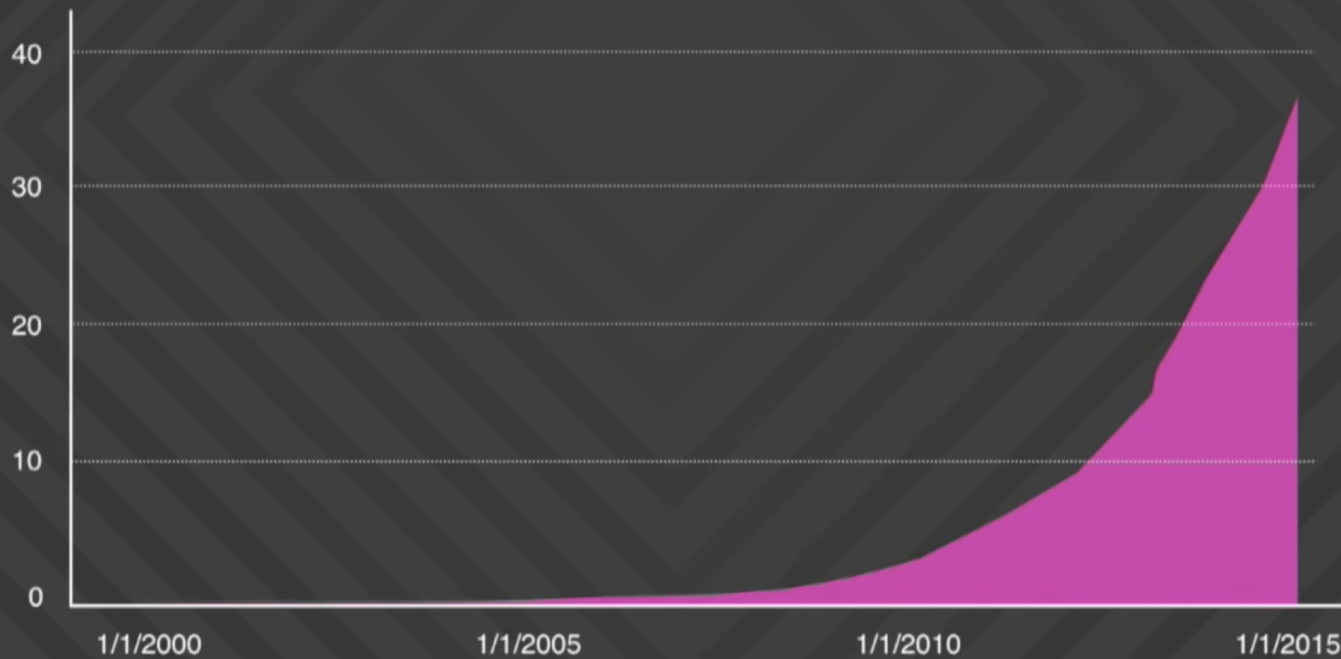
As of Jan 2015

Total number of files*	1 billion
Number of source files	9 million
Lines of code	2 billion
Depth of history	35 million commits
Size of content	86 terabytes
Commits per workday	45 thousand

*The total number of files includes source files copied into release branches, files that are deleted at the latest revision, configuration files, documentation, and supporting data files.

Exponential growth

Millions of changes committed (cumulative)



Google Speed and Scale

- >30,000 developers in 40+ offices
 - 13,000+ projects under active development
 - 30k submissions per day (1 every 3 seconds)
-
- All builds from source
 - 30+ sustained code changes per minute with 90+ peaks
 - 50% of code changes monthly
 - 150+ million test cases / day, > 150 years of test / day
 - Supports continuous deployment for all Google teams!

Google code base vs Linux kernel code base

Some perspective

Linux kernel

- 15 million lines of code in 40 thousand files (total)

Google repository

- 15 million lines of code in 250 thousand files *changed per week, by humans*
- 2 billion lines of code, in 9 million source files (total)

How do they do it?

1. Lots of (automated) testing

Google workflow



- All code is reviewed before commit (by humans and automated tooling)
- Each directory has a set of owners who must approve the change to their area of the repository
- Tests and automated checks are performed before and after commit
- Auto-rollback of a commit may occur in the case of widespread breakage

2. Lots of automation

Additional tooling support

Critique	Code review
CodeSearch*	Code browsing, exploration, understanding, and archeology
Tricorder**	Static analysis of code surfaced in Critique, CodeSearch
Presubmits	Customizable checks, testing, can block commit
TAP	Comprehensive testing before and after commit, auto-rollback
Rosie	Large-scale change distribution and management

* See "How Developers Search for Code: A Case Study", In European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering, 2015

** See "Tricorder: Building a program analysis ecosystem". In International Conference on Software Engineering (ICSE), 2015

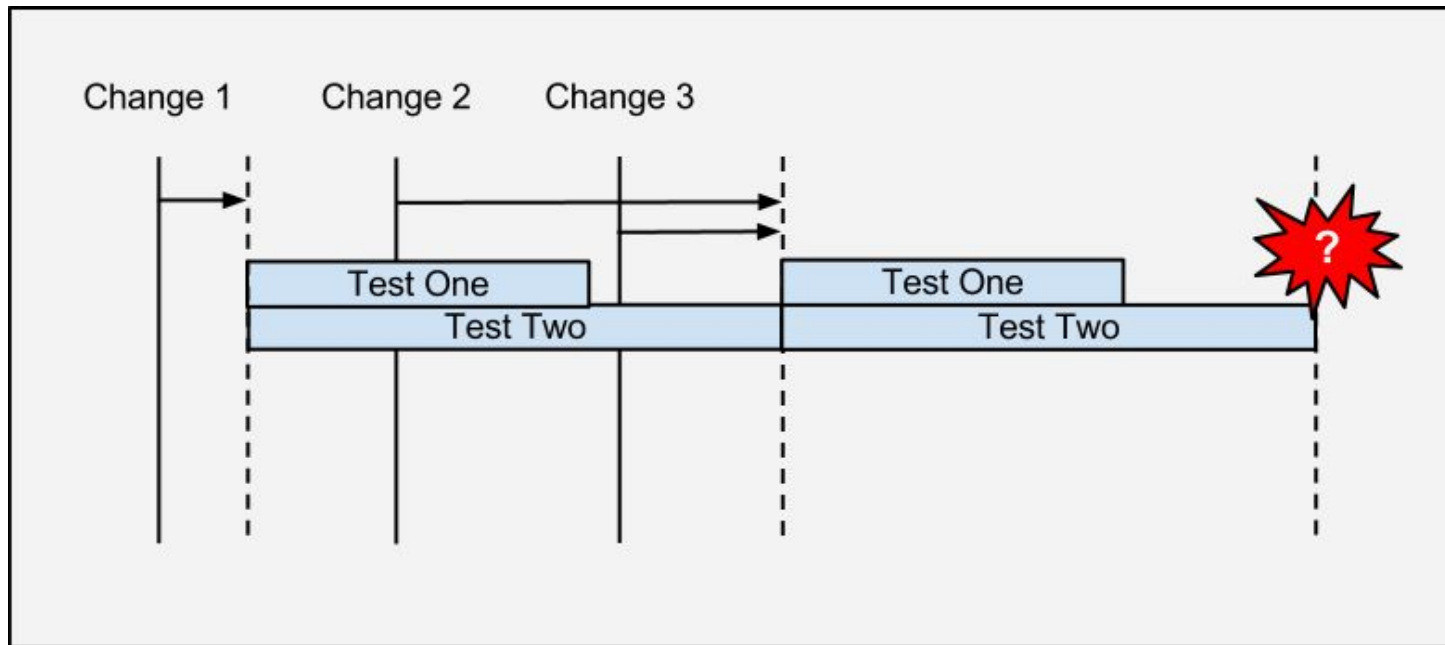
3. Smarter tooling

- Build system
- Version control
- ...

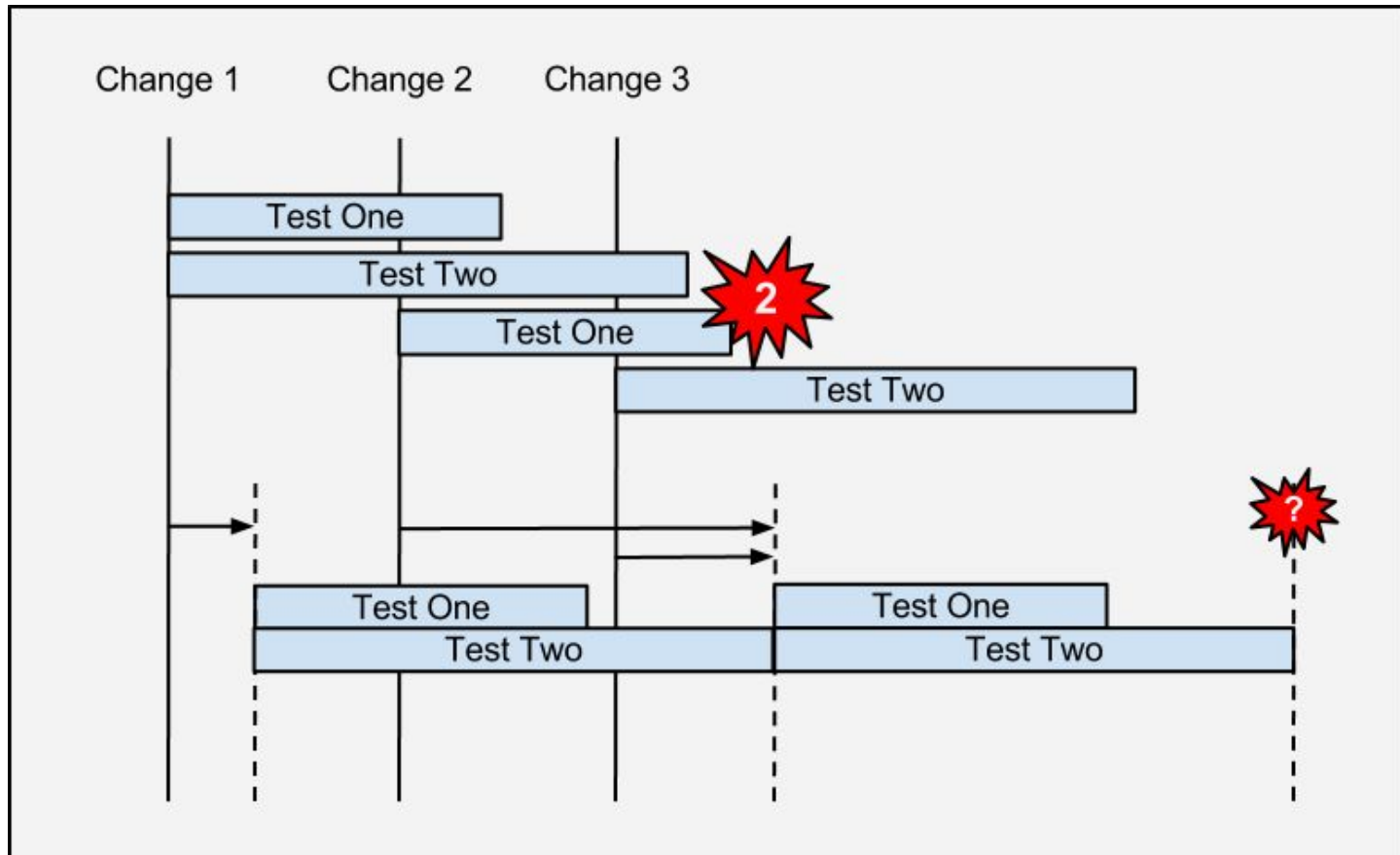
3a. Build system

Google Standard Continuous Build System

- Triggers builds in continuous cycle
- Cycle time = longest build + test cycle
- Tests many changes together
- Which change broke the build?



- Triggers tests on every change
- Uses fine-grained dependencies
- Change 2 broke test 1



Which tests to run?

GMAIL

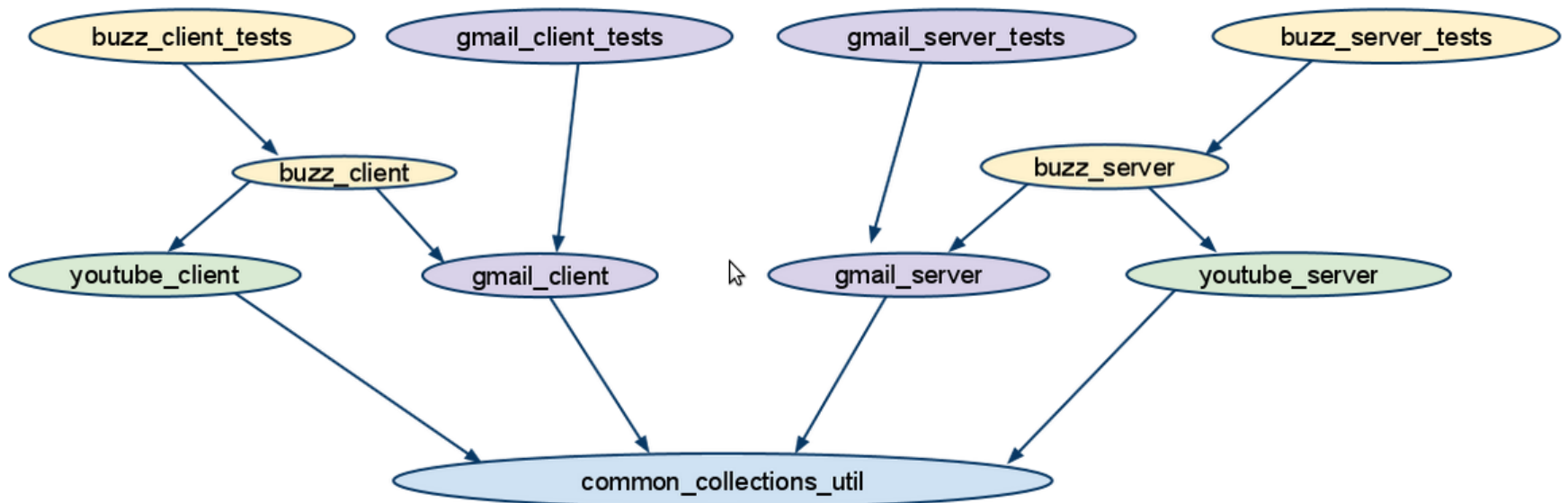
Test Target:

name: //depot/gmail_client_tests
name: //depot/gmail_server_tests

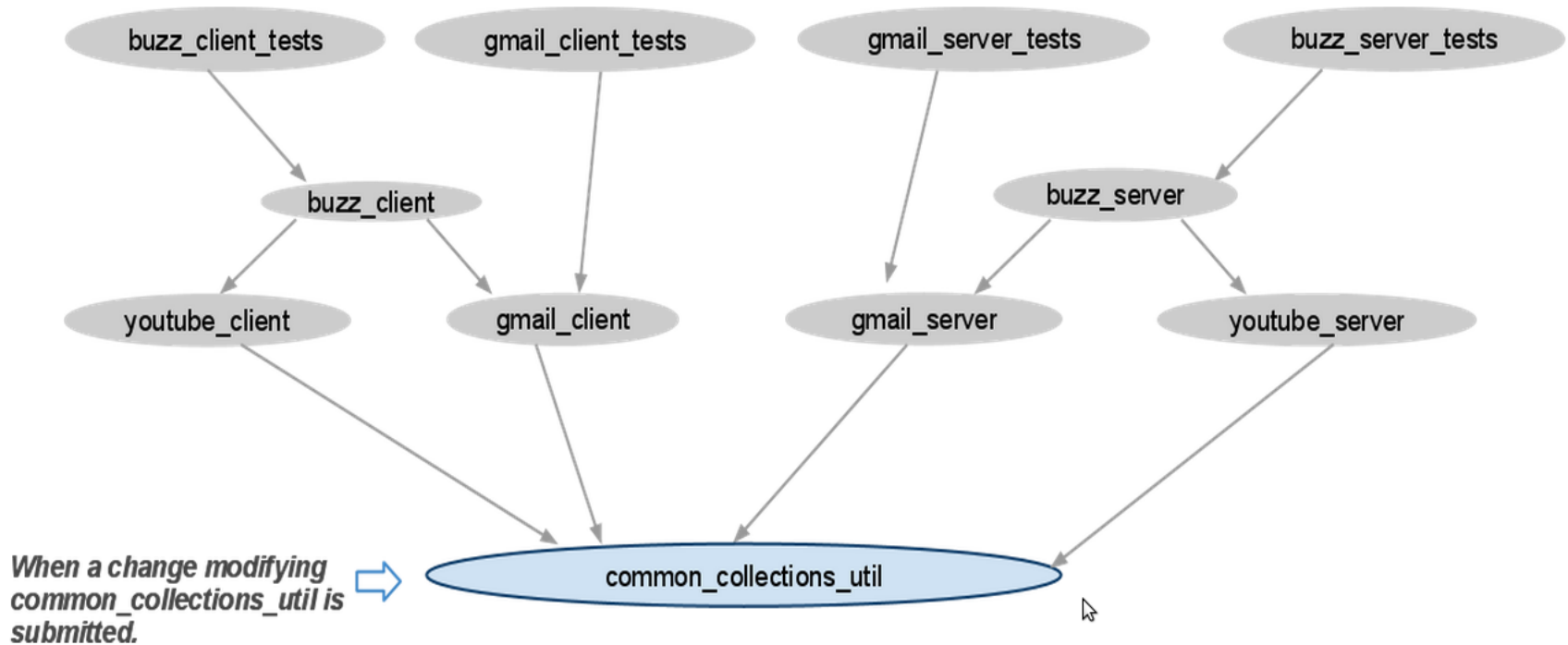
BUZZ

Test targets:

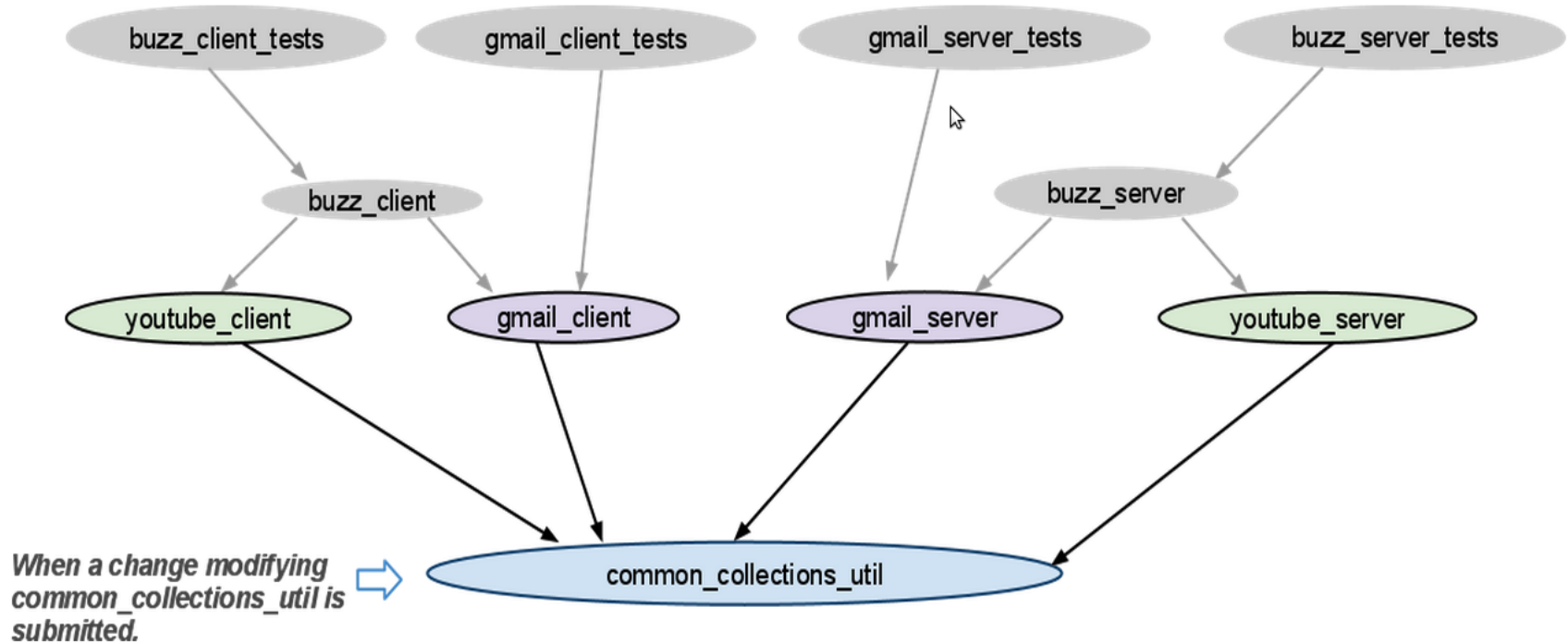
name: //depot/buzz_server_tests
name: //depot/buzz_client_tests



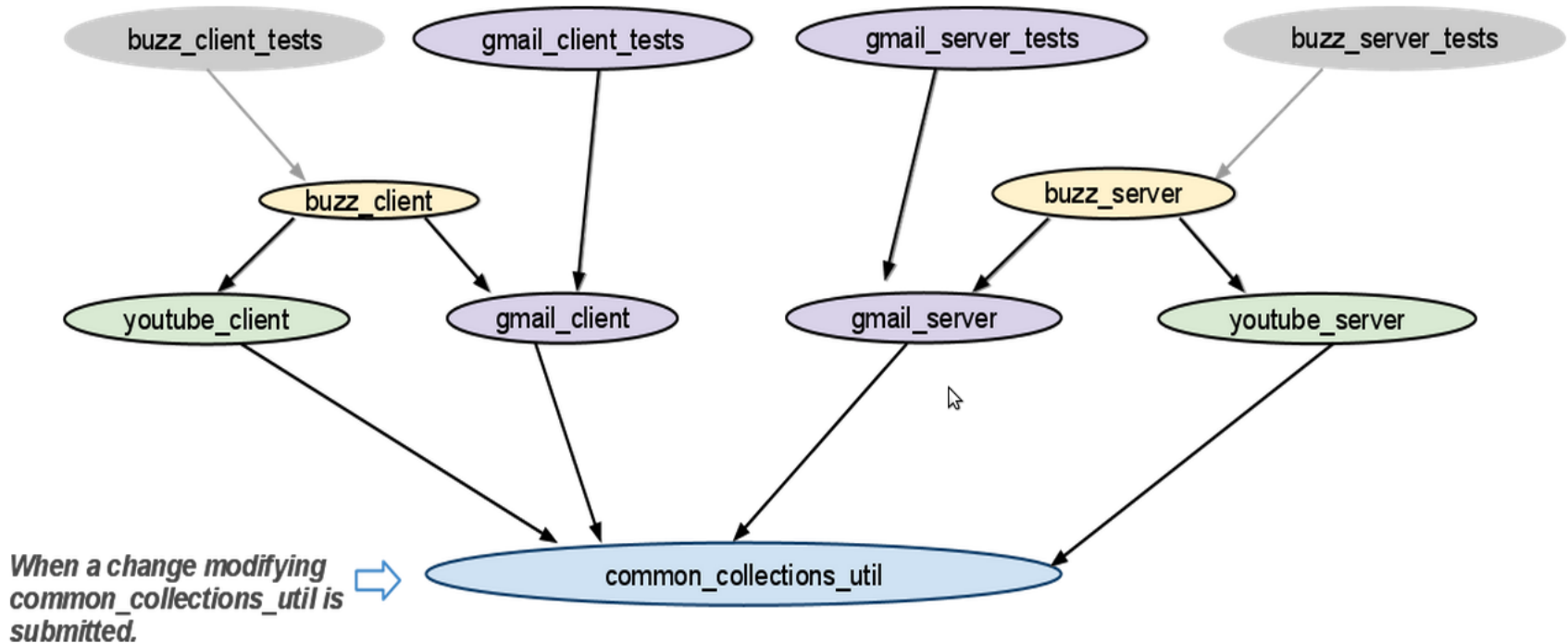
Scenario 1: a change modifies common_collections_util



Scenario 1: a change modifies common_collections_util

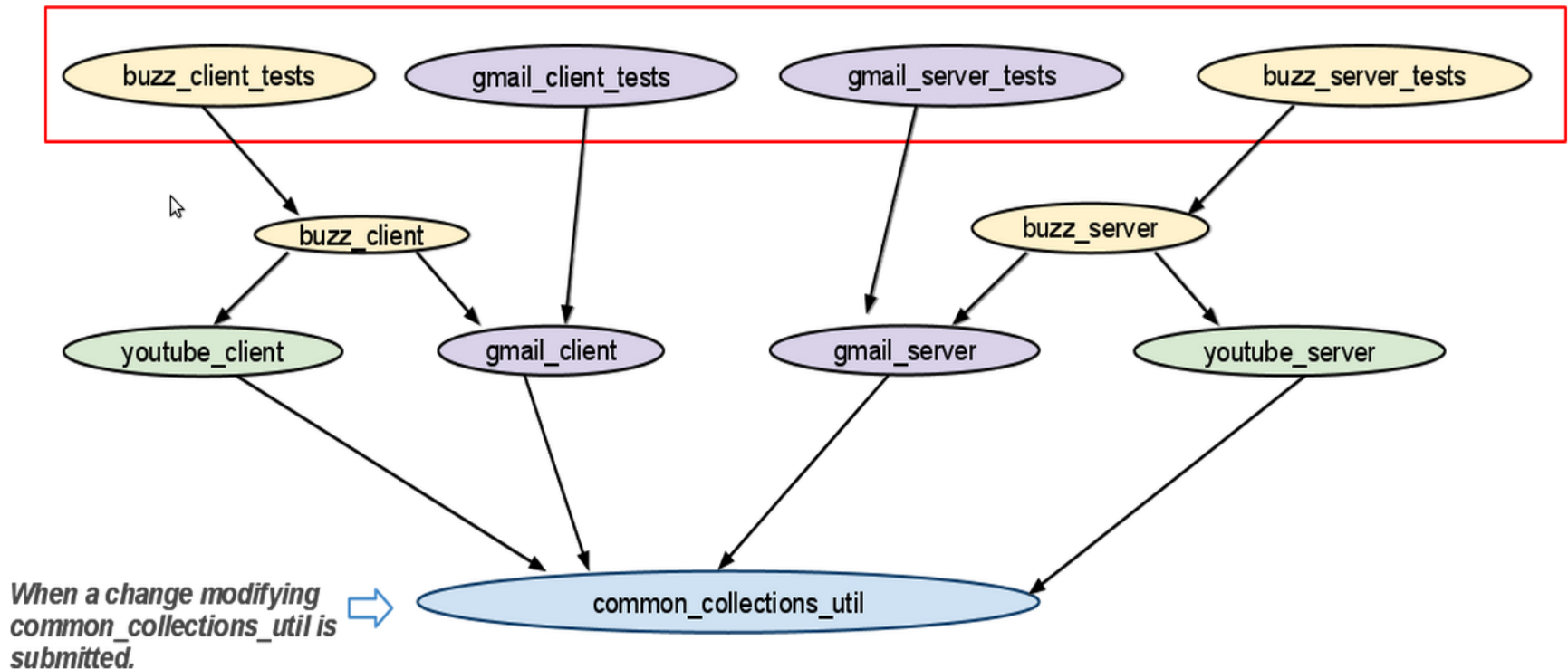


Scenario 1: a change modifies common_collections_util

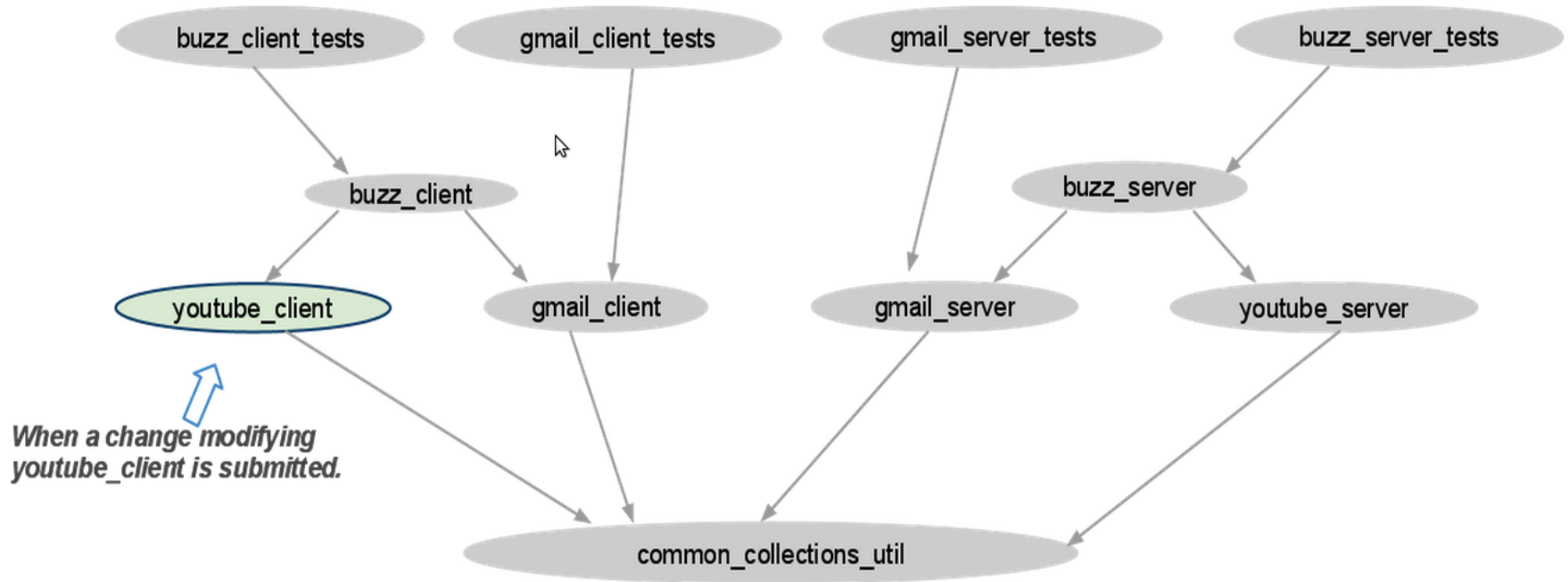


Scenario 1: a change modifies common_collections_util

All tests are affected! Both Gmail and Buzz projects need to be updated

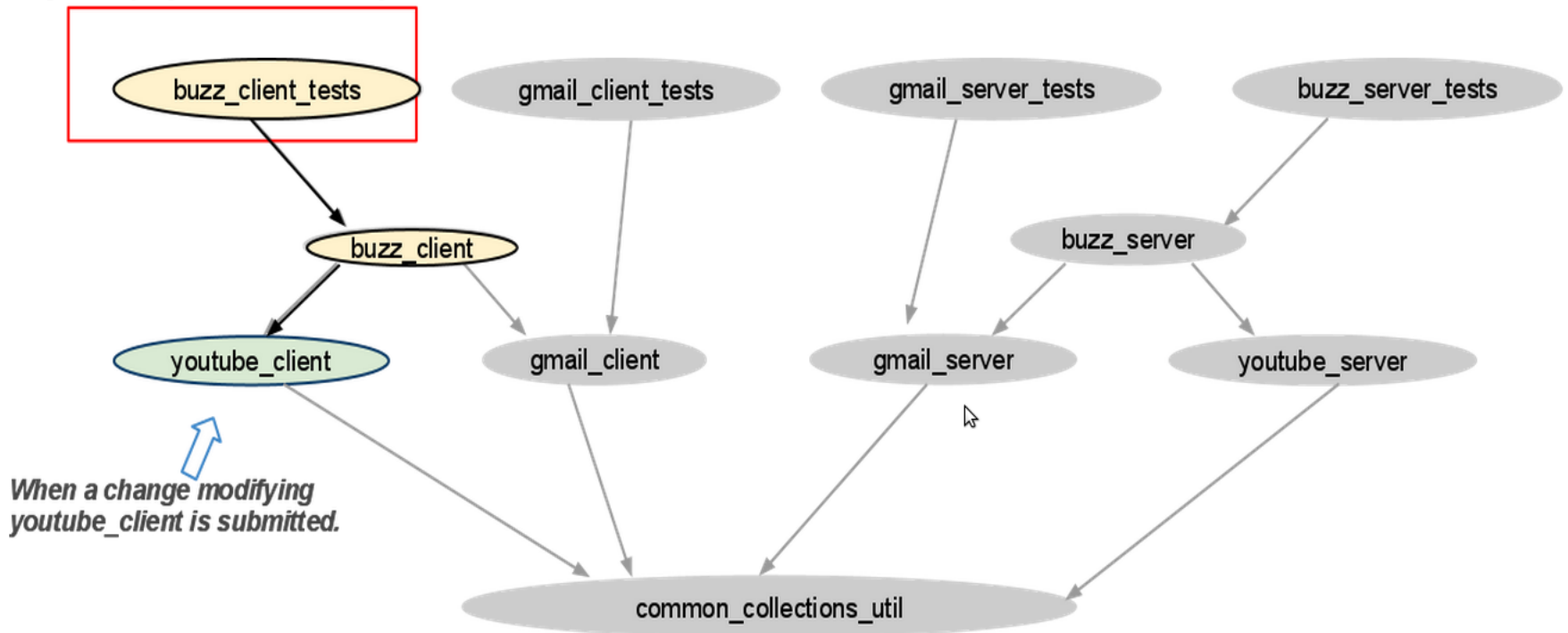


Scenario 2: a change modifies the youtube_client



Scenario 2: a change modifies the youtube_client

Only buzz_client_tests are run and only Buzz project needs to be updated.



3b. Version control

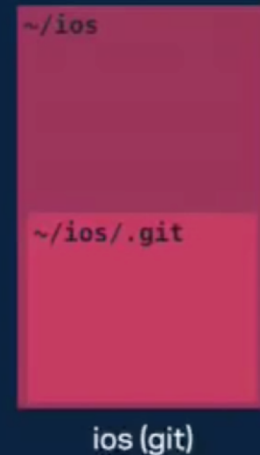
- Problem: even git can get slow at Facebook-like scale
 - 1M+ source control commands run per day
 - 100K+ commits per week

Cloning with git: iOS Today

Many files

Deep history

Large “footprint” makes git slow



3b. Version control

- Solution: redesign version control

Enter Mercurial: Sparse Checkouts

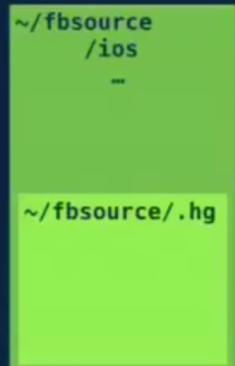
Work on only the files you need.

Build system knows how to
check out more.

Enter Mercurial: Shallow History

Work locally without complete history.

Need more history?
Downloaded automatically on demand.

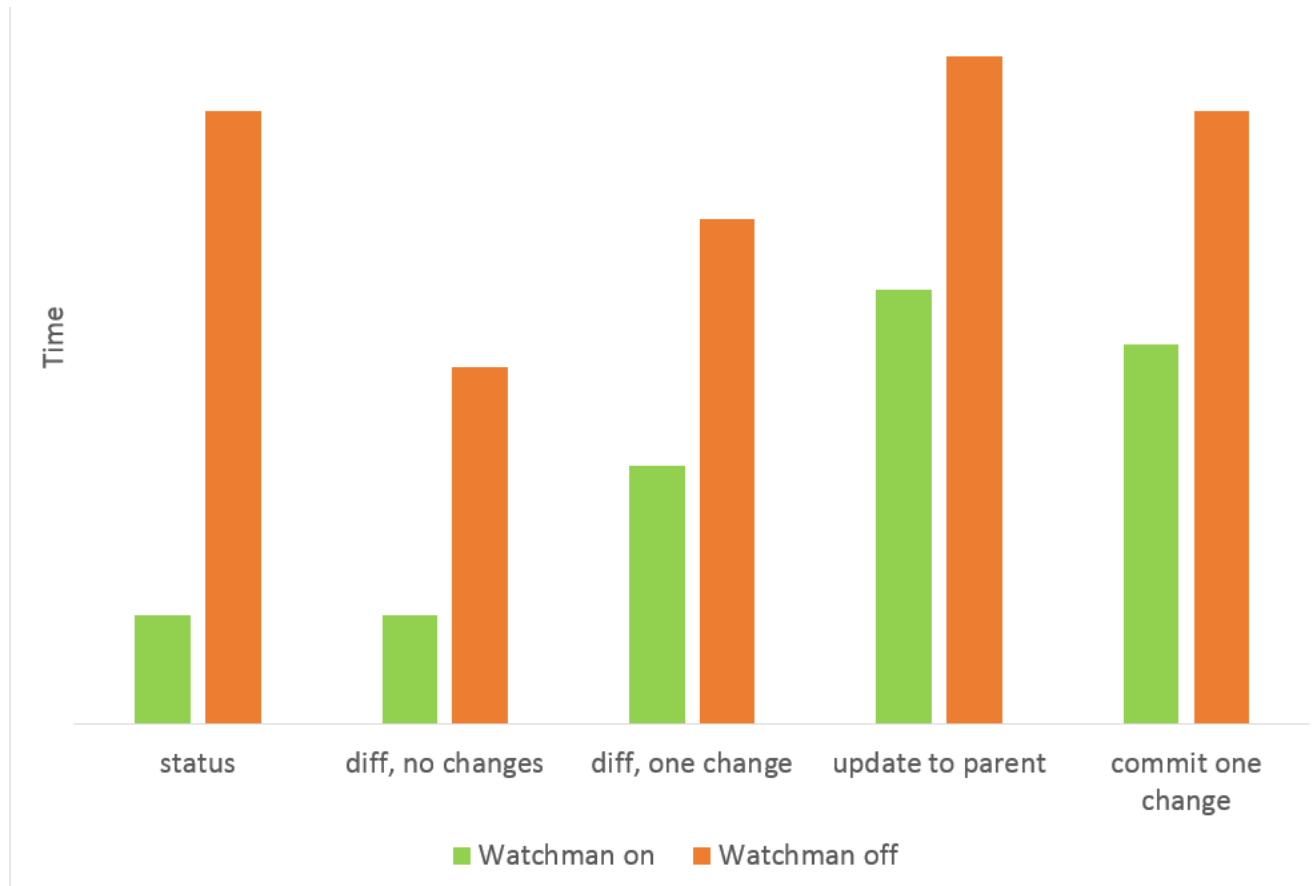


3b. Version control

- Solution: redesign version control
 - Query build system's file monitor, Watchman, to see which files have changed

3b. Version control

- Solution: redesign version control
 - Query build system's file monitor, Watchman, to see which files have changed → **5x faster “status” command**



3b. Version control

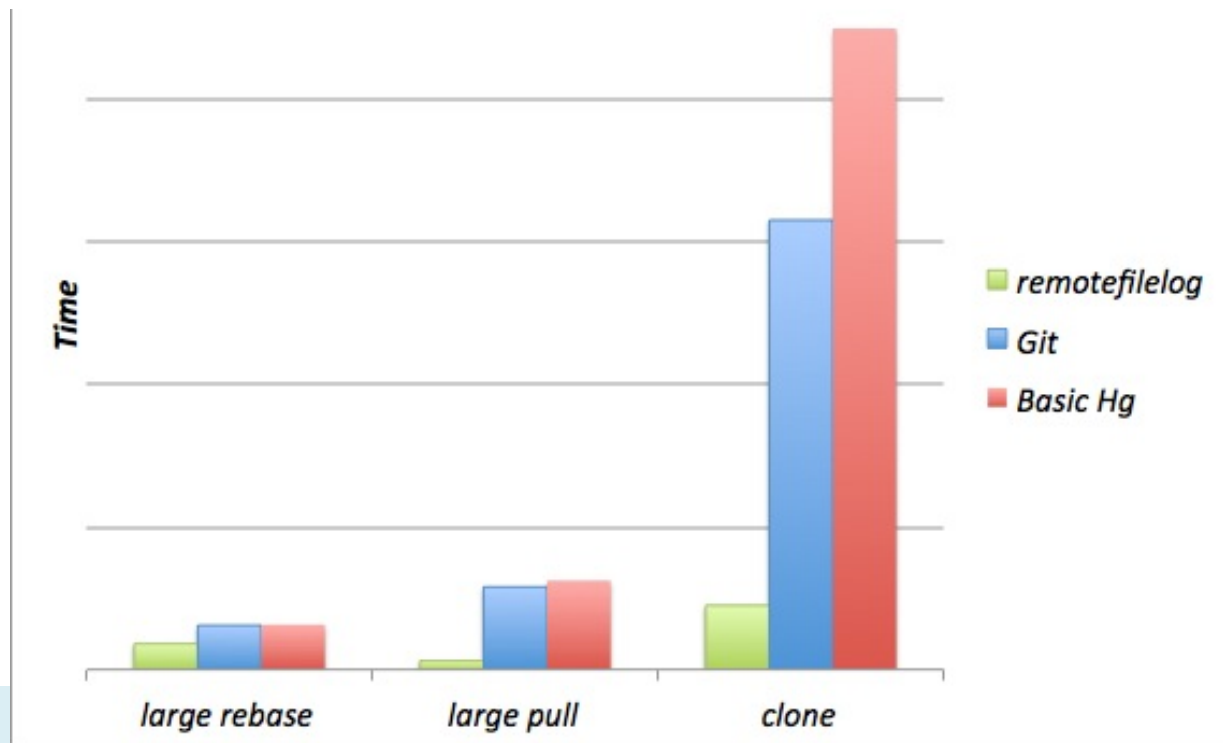
- Solution: redesign version control
 - Sparse checkouts??? (remember, git is a distributed VCS)

3b. Version control

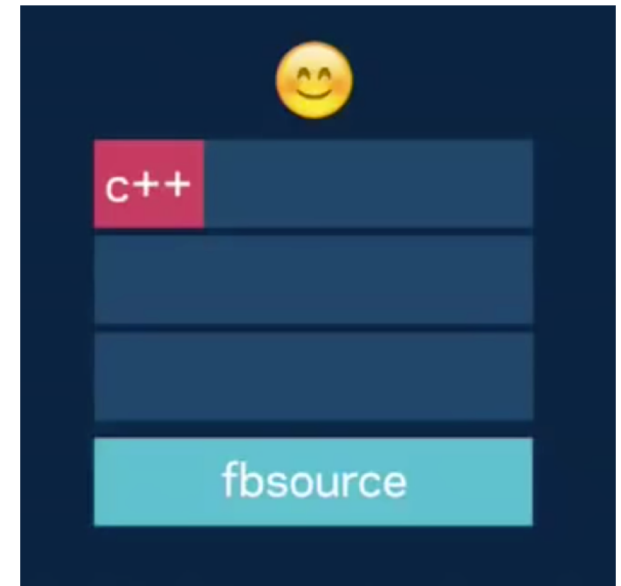
- Solution: redesign version control
 - Sparse checkouts:
 - Change the clone and pull commands to download only the commit metadata, while omitting all file changes (the bulk of the download)
 - When a user performs an operation that needs the contents of files (such as checkout), download the file contents on demand using existing memcache infrastructure

3b. Version control

- Solution: redesign version control
 - Sparse checkouts → **10x faster clones and pulls**
 - Change the clone and pull commands to download only the commit metadata, while omitting all file changes (the bulk of the download)
 - When a user performs an operation that needs the contents of files (such as checkout), download the file contents on demand using existing memcache infrastructure



4. Monolithic repository

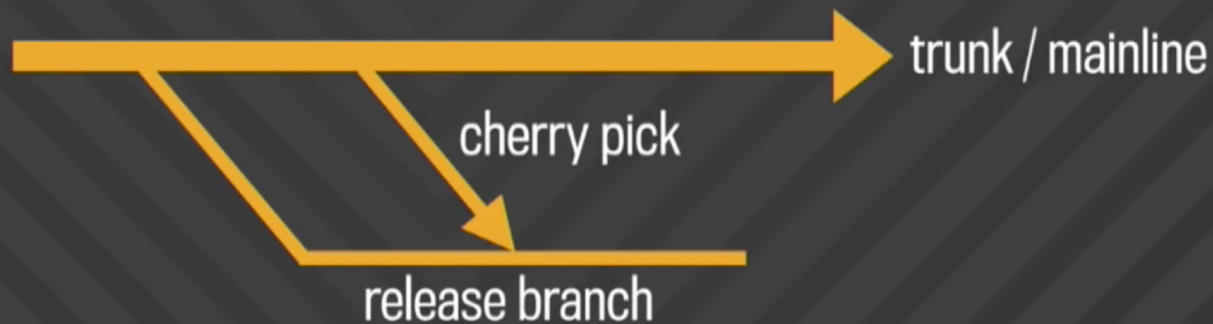


Monolithic repository – no major use of branches for development

Trunk-based development

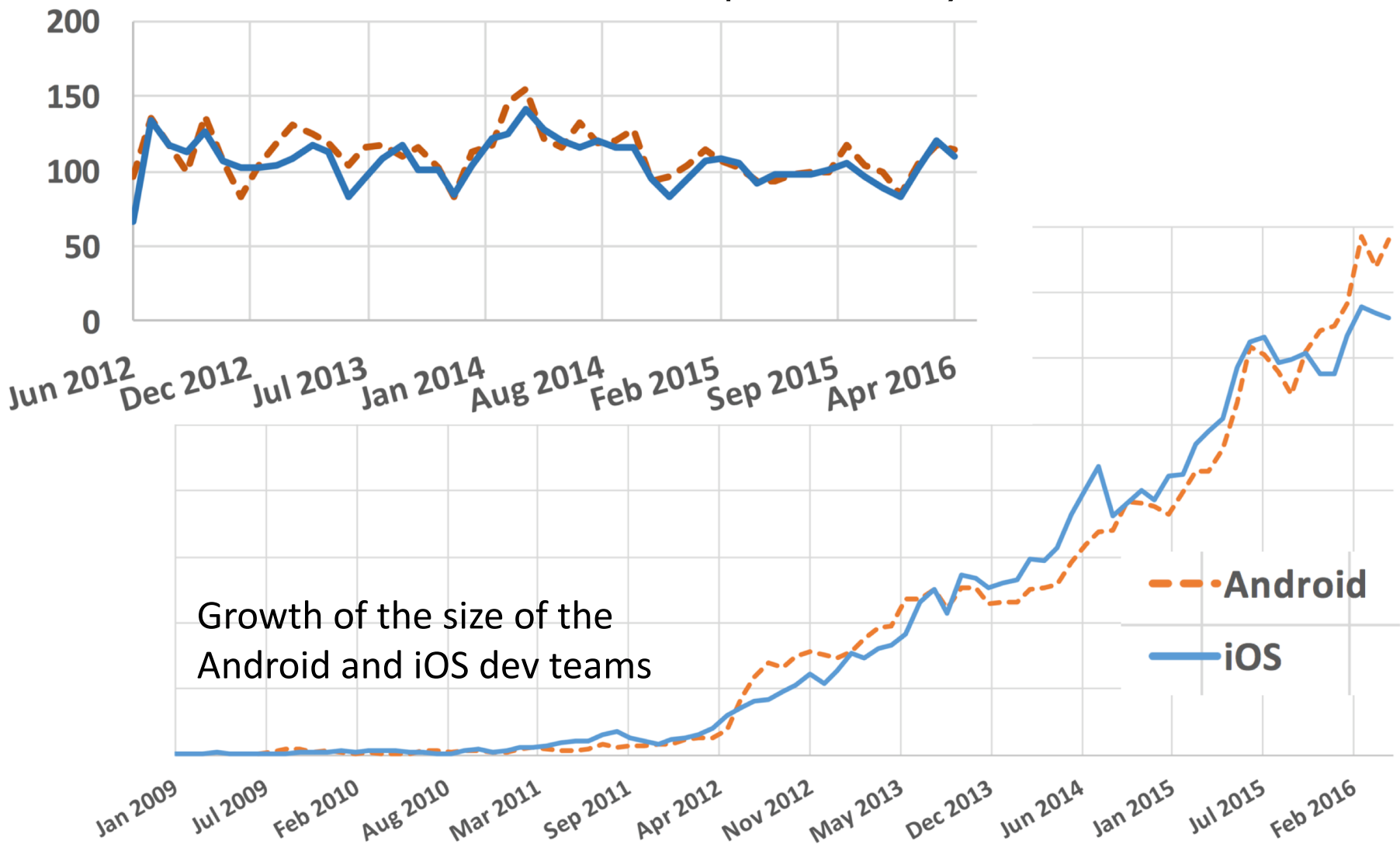
Combined with a centralized repository, this defines the monolithic model

- Piper users work at “head”, a consistent view of the codebase
- All changes are made to the repository in a single, serial ordering
- There is no significant use of branching for development
- Release branches are cut from a specific revision of the repository



Did it work? Yes. Sustained productivity at Facebook

Lines Committed Per Developer Per Day



MONOREPO VS MANY REPOS

A recent history of code organization

- A single team with a monolithic application in a single repository
- ...
- Multiple teams with many separate applications in many separate repositories
- Multiple teams with many ~~separate applications~~ **microservices** in many separate repositories
- A single team with many microservices in many repositories
- ...
- Many teams with many applications in one big **Monorepo**

What is a Monolithic Repository (monorepo)?

A **single** version control repository containing multiple

- ▶ projects
- ▶ applications
- ▶ libraries,


often using a common build system.

Monorepos in industry

Google (computer science version)

TRUSTED INSIGHTS FOR COMPUTING'S LEADING PROFESSIONALS

ACM.org | Join ACM | About Communications | ACM Resources | Alerts & Feeds

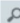
SIGN IN

COMMUNICATIONS

OF THE

ACM

HOME | CURRENT ISSUE | NEWS | BLOGS | OPINION | RESEARCH | PRACTICE | CAREERS | ARCHIVE | VIDEOS

Search 

Home / Magazine Archive / July 2016 (Vol. 59, No. 7) / Why Google Stores Billions of Lines of Code in a Single... / Full Text

CONTRIBUTED ARTICLES

Why Google Stores Billions of Lines of Code in a Single Repository

By Rachel Potvin, Josh Levenberg
Communications of the ACM, Vol. 59 No. 7, Pages 78-87
10.1145/2854146
[Comments \(3\)](#)

VIEW AS:     

SHARE:      



Early Google employees decided to work with a shared codebase managed through a centralized source control system. This approach has served Google well for more than 16 years, and today the vast majority of Google's software assets continues to be stored in a single, shared repository. Meanwhile, the number of Google software developers has steadily increased, and the size of the Google codebase has grown exponentially (see [Figure 1](#)). As a result, the technology used to host the codebase has also evolved significantly.

[Back to Top](#)

[Key Insights](#)

SIGN IN for Full Access

User Name 

Password 

[» Forgot Password?](#)

[» Create an ACM Web Account](#)

SIGN IN

ARTICLE CONTENTS:

[Introduction](#)

[Key Insights](#)

[Google-Scale](#)

[Background](#)

[Analysis](#)

[Alternatives](#)

Advantages and Disadvantages of a Monolithic Repository

A case study at Google

Ciera Jaspan, Matthew Jorde,
Andrea Knight, Caitlin Sadowski,
Edward K. Smith, Collin Winter
Google

ciera,majorde,aknight,supertri,edwardsmith,
collinwinter@google.com

Emerson Murphy-Hill*
NC State University
emerson@csc.ncsu.edu

ABSTRACT

Monolithic source code repositories (repos) are used by several large tech companies, but little is known about their advantages or disadvantages compared to multiple per-project repos. This paper investigates the relative tradeoffs by utilizing a mixed-methods approach. Our primary contribution is a survey of engineers who have experience with both monolithic repos and multiple, per-project repos. This paper also backs up the claims made by these engineers with a large-scale analysis of developer tool logs. Our study finds that the visibility of the codebase is a significant advantage of a monolithic repo: it enables engineers to discover APIs to reuse, find examples for using an API, and automatically have dependent code updated as an API migrates to a new version. Engineers also appreciate the centralization of dependency management in the repo. In contrast, multiple-repository (multi-repo) systems afford engineers more flexibility to select their own toolchains and provide significant access control and stability benefits. In both cases, the related tooling is also a significant factor; engineers favor particular tools and are drawn to repo management systems that support their desired toolchain.

CCS CONCEPTS

• **Software and its engineering** → **Software configuration management and version control systems**;

1 INTRODUCTION

Companies today are producing more source code than ever before. Given the increasingly large codebases involved, it is worth examining the software engineering experience pro-

the organization. Successfully organizing these dependencies and frameworks is crucial for development velocity.

One approach to scaling development practices is the monolithic repo, a model of source code organization where engineers have broad access to source code, a shared set of tooling, and a single set of common dependencies. This standardization and level of access is enabled by having a single, shared repo that stores the source code for all the projects in an organization. Several large software companies have already moved to this organizational model, including Facebook, Google, and Microsoft [10, 12, 17, 21]; however, there is little research addressing the possible advantages or disadvantages of such a model. Does broad access to source code let software engineers better understand APIs and libraries, or overwhelm engineers with use cases that aren't theirs? Do projects benefit from shared dependency versioning, or would engineers prefer more stability for their dependencies? How often do engineers take advantage of the workflows that monolithic repos enable? Do engineers prefer having consistent, shared toolchains or the flexibility of selecting a toolchain for their project?

In this paper, we investigate the experience of engineers working within a monolithic repo and the tradeoffs between using a monolithic repo and a multi-repo codebase. Specifically, this paper seeks to answer two research questions:

- (1) What do developers perceive as the benefits and drawbacks to working in a monolithic versus multi-repo environment?
- (2) To what extent do developers make use of the unique advantages that monolithic repos provide?

Monorepos in industry

Scaling Mercurial at Facebook

The screenshot shows the Facebook Code website. The header includes the Facebook logo, the word 'Code', a search bar, and navigation links for Open Source, Platforms, Infrastructure Systems, Hardware Infrastructure, Video & VR, and Artificial Intelligence. The main content area features a blog post titled 'Scaling Mercurial at Facebook' dated 7 January 2014, with tags for INFRA, OPEN SOURCE, PERFORMANCE, and OPTIMIZATION. The authors are Durham Goode and Siddharth P Agarwal. The post text discusses Facebook's source control challenges and the choice of Mercurial. A 'Recommended' sidebar on the right lists other articles like 'Scaling memcached at Facebook' and 'Flashcache at Facebook: From 2010 to 2013 and beyond'.

f Code Search

Open Source Platforms Infrastructure Systems Hardware Infrastructure Video & VR Artificial Intelligence

7 January 2014 INFRA · OPEN SOURCE · PERFORMANCE · OPTIMIZATION

Scaling Mercurial at Facebook

Durham Goode Siddharth P Agarwal

With thousands of commits a week across hundreds of thousands of files, Facebook's main source repository is enormous—many times larger than even the Linux kernel, which checked in at 17 million lines of code and 44,000 files in 2013. Given our size and complexity—and Facebook's practice of shipping code twice a day—improving our source control is one way we help our engineers move fast.

Choosing a source control system

Two years ago, as we saw our repository continue to grow at a staggering rate, we sat down and extrapolated our growth forward a few years. Based on those projections, it appeared likely that our then-current technology, a Subversion server with a Git mirror, would become a productivity bottleneck very soon. We looked at the available options and found none that were both fast and easy to use at scale.

Our code base has grown organically and its internal dependencies are very complex. We could have spent a lot of time making it more modular in a way that would be friendly to a source control tool, but there are a number of benefits to using a single repository. Even at our current scale, we often make large changes throughout our code base, and having a single repository is useful for continuous

Recommended

- Scaling memcached at Facebook
- Flashcache at Facebook: From 2010 to 2013 and beyond

Monorepos in industry

Microsoft claim the largest git repo on the planet

Server & Tools Blogs > Developer Tools Blogs > Brian Harrys blog

Sign in

Executive Bloggers

Visual Studio

DevOps

Languages

.NET

Platform
Development

Data Development

Brian Harrys blog

Everything you want to know about Visual Studio ALM and Farming

The largest Git repo on the planet

05/24/2017 by Brian Harry MS // 59 Comments

Share 2.2k 3213 1230

It's been 3 months since I first wrote about [our efforts to scale Git to extremely large projects and teams](#) with an effort we called "Git Virtual File System". As a reminder, GVFS, together with a set of enhancements to Git, enables Git to scale to VERY large repos by virtualizing both the .git folder and the working directory. Rather than download the entire repo and checkout all the files, it dynamically downloads only the portions you need based on what you use.

A lot has happened and I wanted to give you an update. Three months ago, GVFS was still a dream. I don't mean it didn't exist – we had a concrete implementation, but rather, it was unproven. We had validated on some big repos but we hadn't rolled it out to any meaningful number of engineers so we had only conviction that it was going to work. Now we have proof.

Today, I want to share our results. In addition, we're announcing the next steps in our GVFS journey for customers, including expanded open sourcing to start taking contributions and improving how it works for us at Microsoft, as well as for partners and customers.

Windows is live on Git

Over the past 3 months, we have largely completed the rollout of Git/GVFS to the Windows team at Microsoft.

As a refresher, the Windows code base is approximately 3.5M files and, when checked in to a Git repo, results in a repo of about 300GB.

Visual Studio

Download Visual Studio →

Download TFS →

Visual Studio Team Services →

Search

Search MSDN with Bing

☐ Search this blog ☒ Search all blogs

Subscribe Blog via Email

Subscribe to this blog and receive notifications of new posts by email.

Email Address

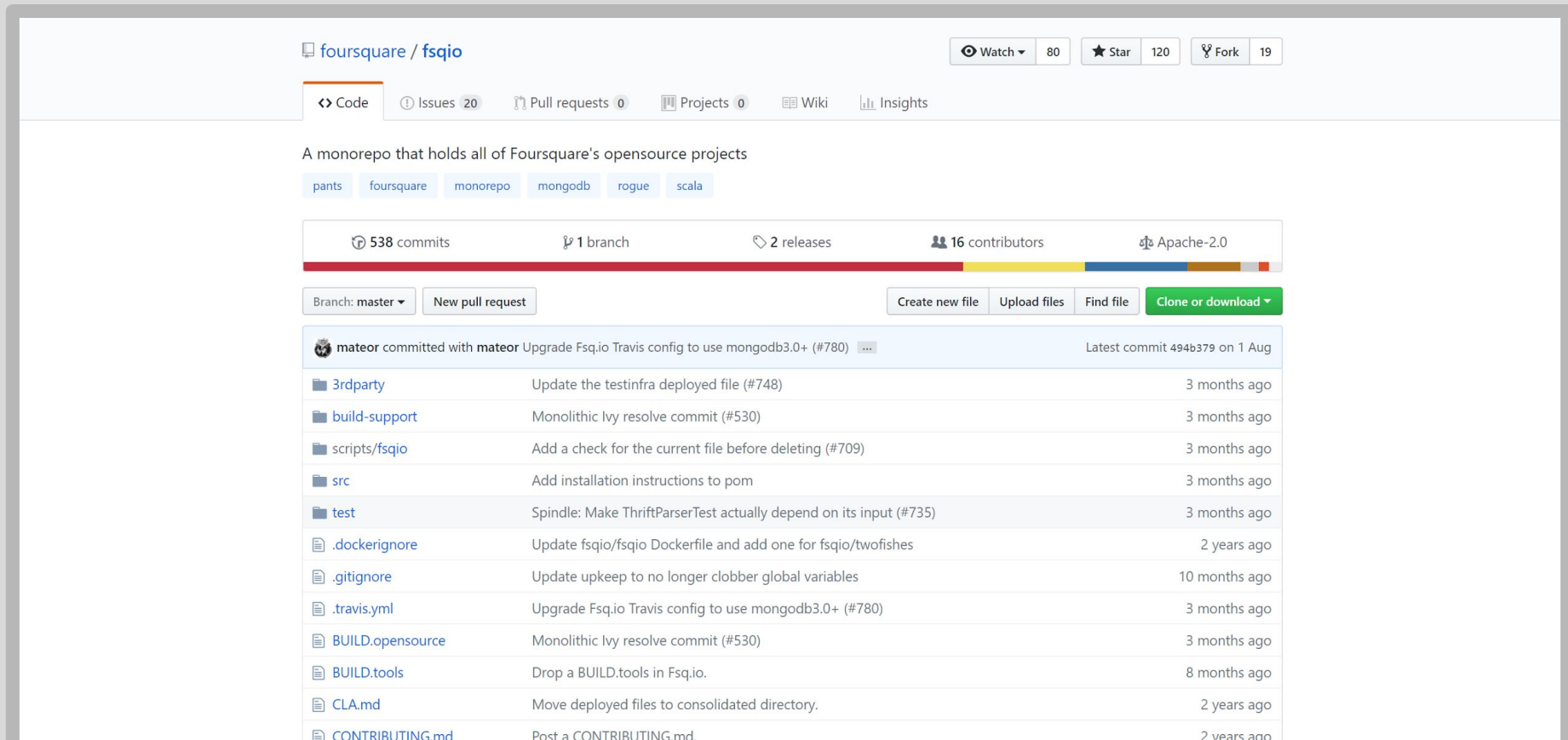
Subscribe!

Unsubscribe



Monorepos in open-source

foresquare public monorepo



The screenshot displays the GitHub repository page for `foursquare/fsqio`. At the top, the repository name is shown with navigation links for Code, Issues (20), Pull requests (0), Projects (0), Wiki, and Insights. On the right, there are buttons for Watch (80), Star (120), and Fork (19). Below the repository name, a description states: "A monorepo that holds all of Foursquare's opensource projects". A list of sub-projects is provided: pants, foursquare, monorepo, mongodb, rogue, and scala. A progress bar shows the repository's activity: 538 commits, 1 branch, 2 releases, 16 contributors, and Apache-2.0 license. Below the progress bar, there are buttons for "Branch: master", "New pull request", "Create new file", "Upload files", "Find file", and "Clone or download". The commit history is listed below, showing the latest commit by `mateor` on August 1st, 2016, with the message "Upgrade Fsquito Travis config to use mongodb3.0+ (#780)". The commit history table includes the following entries:

File	Commit Message	Time Ago
3rdparty	Update the testinfra deployed file (#748)	3 months ago
build-support	Monolithic Ivy resolve commit (#530)	3 months ago
scripts/fsqio	Add a check for the current file before deleting (#709)	3 months ago
src	Add installation instructions to pom	3 months ago
test	Spindle: Make ThriftParserTest actually depend on its input (#735)	3 months ago
.dockerignore	Update fsqio/fsqio Dockerfile and add one for fsqio/twofishes	2 years ago
.gitignore	Update upkeep to no longer clobber global variables	10 months ago
.travis.yml	Upgrade Fsquito Travis config to use mongodb3.0+ (#780)	3 months ago
BUILD.opensource	Monolithic Ivy resolve commit (#530)	3 months ago
BUILD.tools	Drop a BUILD.tools in Fsquito.	8 months ago
CLA.md	Move deployed files to consolidated directory.	2 years ago
CONTRIBUTING.md	Post a CONTRIBUTING.md.	2 years ago

Monorepos in open-source

The Symfony monorepo

43 projects, **25 000** commits, and **400 000** LOC

<https://github.com/symfony/symfony>

Bridge/

5 sub-projects

Bundle/

5 sub-projects

Component/

33 independent sub-projects like Asset, Cache, CssSelector, Finder, Form, HttpKernel, Ldap, Routing, Security, Serializer, Templating, Translation, Yaml, ...

Common build system

Bazel from Google

[Documentation](#)[Contribute](#)[Blog](#)[GitHub](#)

Buck from Facebook

GET I

Pants from Twitter

[Pants](#)[Docs](#)[Community](#)[GitHub](#)

Getting Started

[Installing Pants](#)[Setting Up Pants](#)[Tutorial](#)[Common Tasks](#)

Pants Basics

[Why Use Pants?](#)[Pants Concepts](#)[BUILD files](#)[Target Addresses](#)[Third-Party Dependencies](#)[Pants Options](#)[Invoking Pants](#)[Reporting Server](#)[IDE Support](#)

JVM

[JVM Projects with Pants](#)[JVM 3rdparty Pattern](#)[Scala Support](#)[Publishing Artifacts](#)[Pants for Maven Experts](#)

Pants: A fast, scalable build system

Pants is a build system designed for codebases that:

- Are large and/or growing rapidly.
- Consist of many subprojects that share a significant amount of code.
- Have complex dependencies on third-party libraries.
- Use a variety of languages, code generators and frameworks.

Pants supports Java, Scala, Python, C/C++, Go, Javascript/Node, Thrift, Protobuf and Android code. Adding support for other languages, frameworks and code generators is straightforward.

Pants is a collaborative open-source project, built and used by Twitter, Foursquare, Square, Medium and [other companies](#).



Getting Started

- [Installing Pants](#)
- [Setting Up Pants](#)
- [Tutorial](#)

Cookbook

The [Common Tasks](#) documentation is a practical, solutions-oriented guide to some of the Pants tasks that you're most likely to carry out on a daily basis.

Some advantages of monorepos

High Discoverability For Developers

- ▶ Developers can read and explore the whole codebase
- ▶ grep, IDEs and other tools can search the whole codebase
- ▶ IDEs can offer auto-completion for the whole codebase
- ▶ Code Browsers can links between all artifacts in the codebase

Code-Reuse is cheap

Almost zero cost in introducing a new library

- ▶ Extract library code into a new directory/component
- ▶ Use library in other components
- ▶ Profit!

Refactorings in one commit

Allow large scale refactorings with one single, atomic, history-preserving commit

- ▶ Extract Library/Component
- ▶ Rename Functions/Methods/Components
- ▶ Housekeeping (phpcs-fixer, Namespacing, ...)

Another refactoring example

- Make large backward incompatible changes easily... especially if they span different parts of the project
- For example, old APIs can be removed with confidence
 - Change an API endpoint code **and** all its usages in **all** projects in **one** pull request

Some more advantages

- Easy continuous integration and code review for changes spanning several projects
- (Internal) dependency management is a non-issue
- Less context switching for developers
- Code more reusable in other contexts
- Access control is easy

Some downsides

- Require collective responsibility for team and developers
- Require trunk-based development
 - Feature toggles are technical debt (recall financial services example)
- Force you to have only one version of everything
- Scalability requirements for the repository
- Can be hard to deal with updates around things like security issues
- Build and test bloat without very smart build system
- Slow VCS without very smart system
- Permissions?

Summary

- Software development at scale requires a lot of infrastructure
 - Version control, build managers, testing, continuous integration, deployment, ...
- It's hard to scale development
 - Move towards heavy automation (DevOps)
- Continuous deployment increasingly common
- Opportunities from quick release, testing in production, quick rollback