

Reproducible Data Science

Data 1030. Brown University

Isabel Restrepo. Lead Data Scientist. CCV Brown University. October 18th 2020

Slide credits: August Guang, Bradford Roarr, Paul Stey and Andrew Leith

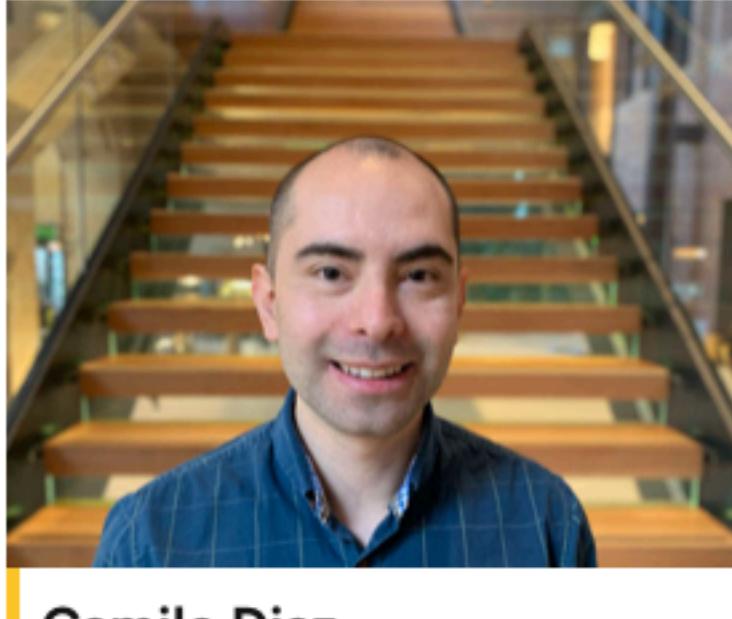
My Team



Benjamin Knörlein
Senior Research Software Engineer
Data & Visualization



Bradford N. Roarr
Senior Research Software Engineer
Data & Visualization



Camilo Diaz
Graphics/Visualization Software Engineer
Data & Visualization



Fernando Gelin
Senior Research Software Engineer
Data & Visualization



Mary McGrath
Senior Research Software Engineer
Data & Visualization



Maura Driscoll
Web Development Intern
Data & Visualization



Rashi Dhar
Web Development Intern
Data & Visualization



What is Reproducible Data Science?

Three main components

- Reproducible Software
- Reproducible Data
- Reproducible Workflows

Dilemma

- Scientific method presumes replicability
- Replication studies have historically NOT been viewed as innovative and often rejected as incremental

Contributing factors

- Competition for tenure track positions
- “Publish or perish”
- Accumulating publications is critical, particularly early in career
- Academic works have no plan for long term maintenance

Reproducibility is essential for replication

- “Methods” section in journal articles
- Historically, this might have satisfied many concerns
- Without technical infrastructure, replication cannot be complete

Efforts to do something about it

- “Why Most Published Research Findings Are False”
 - J. Ioannidis (2005) PLoS Medicine
 - Disregard for statistical rigor, (“*P-*hacking”, Test many hypotheses, no Type-I error correction, small sample, low-power studies)
- FAIR guiding principles (Findable, Accessible, Interoperable and Reusable)
 - NIH
 - Push for papers backed by repositories of software and data.
 - Trends for "live" paper

Reproducible Software

- Version Control
- Ease of use
 - Readable
 - Modular
 - Documentation
- Rigor
 - Testing
 - CI
- Packaging code
 - Cross platform
 - Containerization

Version Control

"A component of software configuration management, also known as revision control or source control, is the management of changes to documents, computer programs, large web sites, and other collections of information."

Wikipedia

Version Control

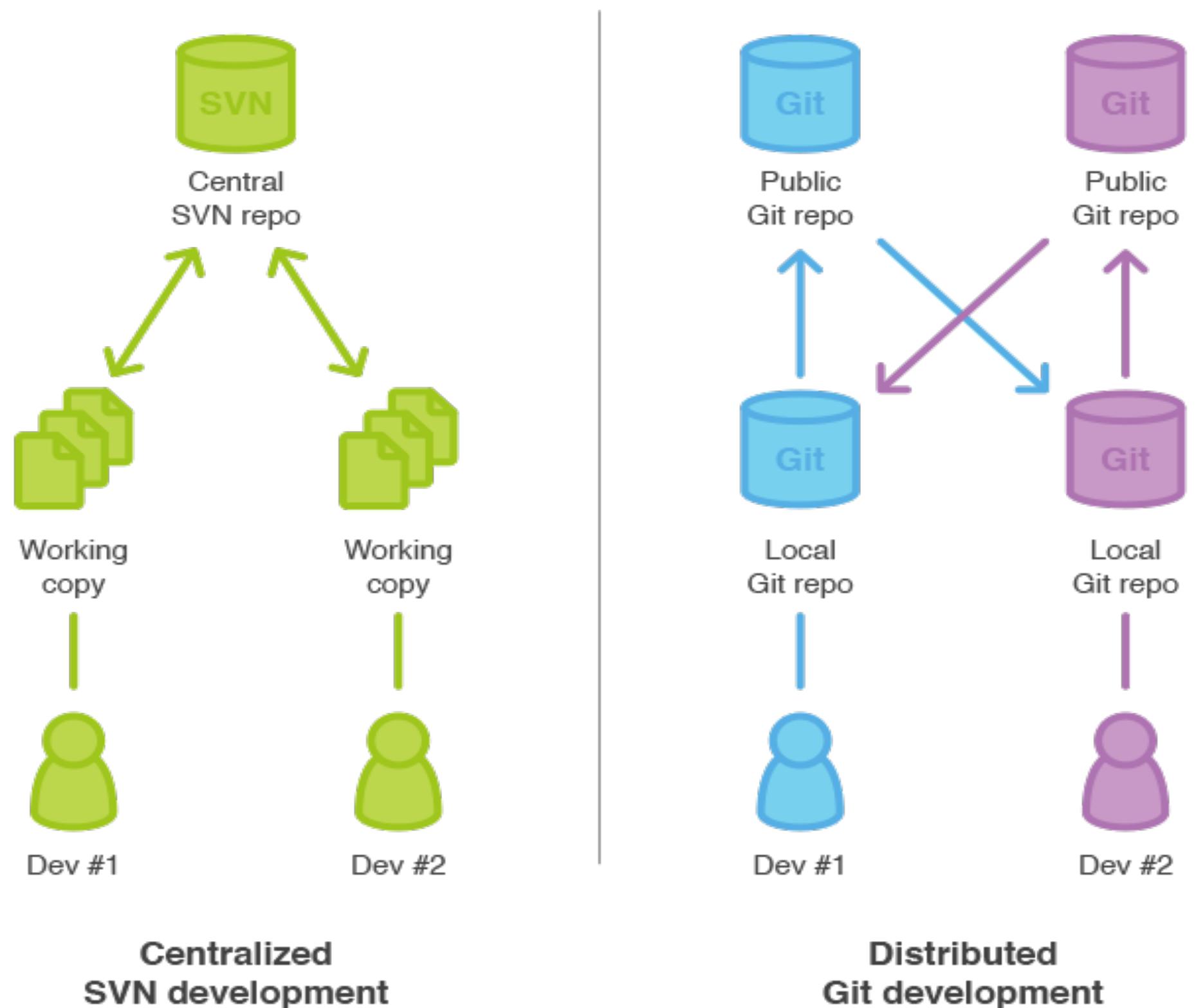
Any software related project, independent of size, should be under version control or in a *repository*

WHY?

- Tracking changes over time
- Good collaboration practices
- Protect stable/production code from bugs

● Technologies

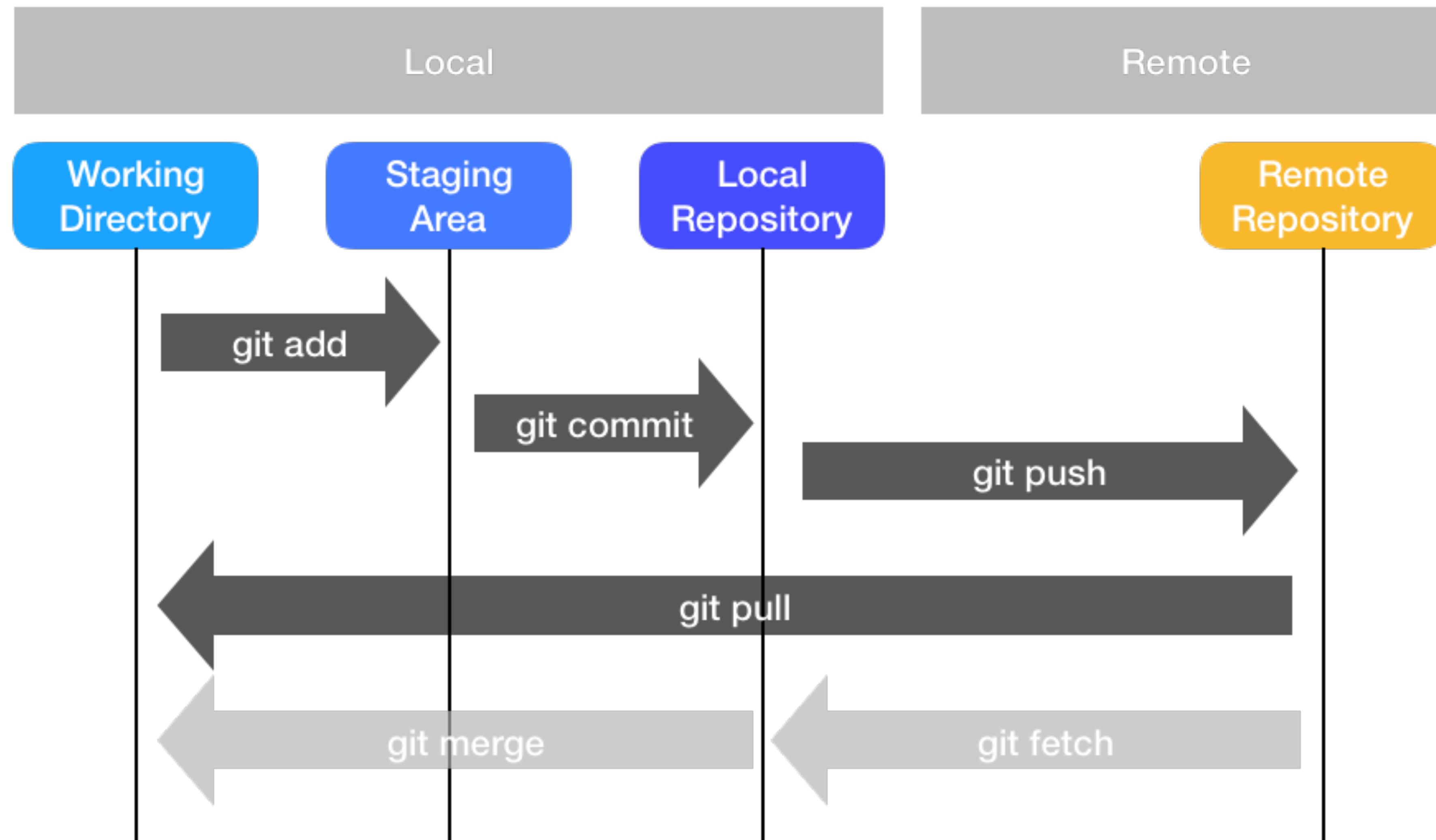
- ~~CVS, SVN, Mercurial (History)~~
- **GIT**



● Hosts



Git Basics



Popular Git Workflows



Feature Workflow

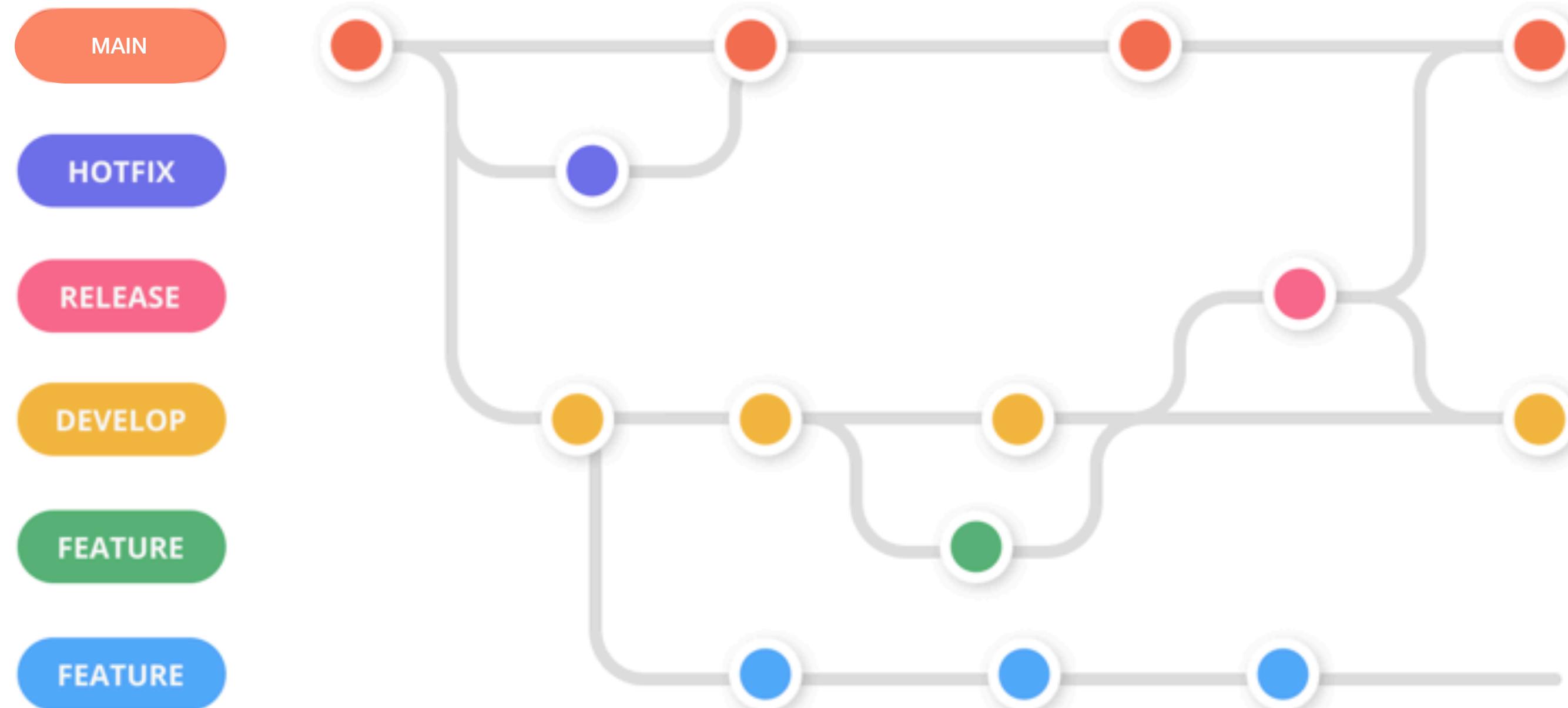
- Create feature/topic branches of main branch
- Main branch is always stable



Feature Workflow + Develop

- Create feature/topic branches of develop branch
- Main branch is the production branch
- Useful for staging/production services

Popular Git Workflows



Gitflow

- Adds hotfix and release branches
- Original workflow
- Popularity has decreased due to complexity

Git Social

 Watch ▾

5,795

 Star

94,463

 Fork

17,812

- Today Version Control is more than simply tracking changes
- It is social:
 - Accelerates rate, quality and number of people working on Open Source
 - Less wheels reinvented :)
 - A way to build your life CV of code
 - Development of tools for testing, documenting, etc that are better integrated

Reproducible Software

- Version Control
- **Ease of use**
 - Readable
 - Modular
 - Documentation
- Rigor
 - Testing
 - CI
- Packaging code
 - Cross platform
 - Containerization

What is READABLE code?

Use a linter: flake8

Use a formatter: black

Use meaningful names for variables

READABLE code

Use comments

Remove commented out code

Simplify statements

Don't abuse one liners

Use Types if possible:
mypy

Modular code

Use small functions

Think of support for
multiple types

Familiarize yourself
with best practices
for specific language.
Functional vs Object
Oriented

Documentation

- ALWAYS include a README.md
 - Basic description
 - Installation
 - Minimum Example
- For more comprehensive documentation use a static-site generator. Host it with repository, i.e. GitHub Pages
 - Track and easy to publish

Documentation

- Basic description
- Installation
- Minimum Example

...

- Badges

- Contributing

Guidelines

README.md



TensorFlow

Documentation	Linux CPU	Linux GPU	Mac OS CPU	Windows CPU	Android
api reference	build passing	build passing	build passing	build failing	build failing Download 1.8.0-rc0

TensorFlow is an open source software library for numerical computation using data flow graphs. The graph nodes represent mathematical operations, while the graph edges represent the multidimensional data arrays (tensors) that flow between them. This flexible architecture enables you to deploy computation to one or more CPUs or GPUs in a desktop, server, or mobile device without rewriting code. TensorFlow also includes [TensorBoard](#), a data visualization toolkit.

TensorFlow was originally developed by researchers and engineers working on the Google Brain team within Google's Machine Intelligence Research organization for the purposes of conducting machine learning and deep neural networks research. The system is general enough to be applicable in a wide variety of other domains, as well.

Keep up to date with release announcements and security updates by subscribing to announce@tensorflow.org.

Documentation

- Basic description
- Installation
- Minimum Example
- ...
- Badges
- Contributing Guidelines

README.md

Installation

See [Installing TensorFlow](#) for instructions on how to install our release binaries or how to build from source.

People who are a little more adventurous can also try our nightly binaries:

Nightly pip packages

- We are pleased to announce that TensorFlow now offers nightly pip packages under the [tf-nightly](#) and [tf-nightly-gpu](#) project on pypi. Simply run `pip install tf-nightly` or `pip install tf-nightly-gpu` in a clean environment to install the nightly TensorFlow build. We support CPU and GPU packages on Linux, Mac, and Windows.

Individual whl files

- Linux CPU-only: [Python 2 \(build history\)](#) / [Python 3.4 \(build history\)](#) / [Python 3.5 \(build history\)](#) / [Python 3.6 \(build history\)](#)
- Linux GPU: [Python 2 \(build history\)](#) / [Python 3.4 \(build history\)](#) / [Python 3.5 \(build history\)](#) / [Python 3.6 \(build history\)](#)
- Mac CPU-only: [Python 2 \(build history\)](#) / [Python 3 \(build history\)](#)
- Windows CPU-only: [Python 3.5 64-bit \(build history\)](#) / [Python 3.6 64-bit \(build history\)](#)
- Windows GPU: [Python 3.5 64-bit \(build history\)](#) / [Python 3.6 64-bit \(build history\)](#)
- Android: [demo APK](#), [native libs \(build history\)](#)

Documentation

- Basic description
- Installation
- Minimum Example

• • •

- Badges
- Contributing Guidelines

README.md

Try your first TensorFlow program

```
$ python
```

```
>>> import tensorflow as tf
>>> hello = tf.constant('Hello, TensorFlow!')
>>> sess = tf.Session()
>>> sess.run(hello)
'Hello, TensorFlow!'
>>> a = tf.constant(10)
>>> b = tf.constant(32)
>>> sess.run(a + b)
42
>>> sess.close()
```

README.md

Contribution guidelines

If you want to contribute to TensorFlow, be sure to review the [contribution guidelines](#). This project adheres to TensorFlow's [code of conduct](#). By participating, you are expected to uphold this code.

We use [GitHub issues](#) for tracking requests and bugs. So please see [TensorFlow Discuss](#) for general questions and discussion, and please direct specific questions to [Stack Overflow](#).

The TensorFlow project strives to abide by generally accepted best practices in open-source software development:

cii best practices in progress 98%

Static-site generators



SLATE
API DOCS GENERATOR



Docusaurus



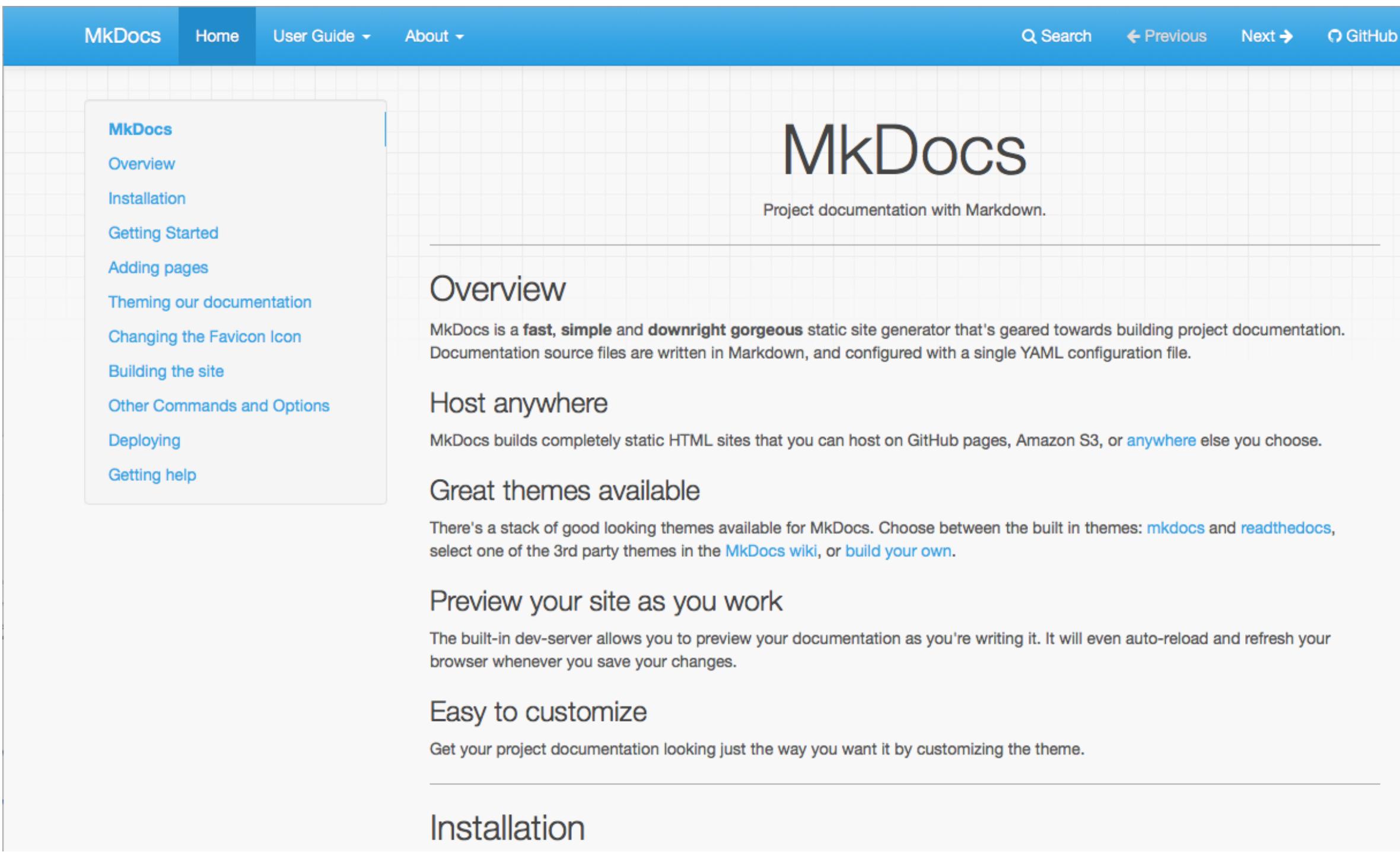
Read the Docs



VuePress

Static-site generators

- From Markdown to beautiful websites out of the box



The screenshot shows the MkDocs documentation site. At the top, there's a navigation bar with links for MkDocs, Home, User Guide, and About. Below the navigation is a search bar and links for Previous and Next. The main content area features the MkDocs logo and the tagline "Project documentation with Markdown." It includes sections for Overview, Host anywhere, Great themes available, Preview your site as you work, Easy to customize, and Installation.

- MkDocs
- Home
- User Guide ▾
- About ▾

Q Search ← Previous Next → GitHub

MkDocs

Project documentation with Markdown.

Overview

MkDocs is a fast, simple and downright gorgeous static site generator that's geared towards building project documentation. Documentation source files are written in Markdown, and configured with a single YAML configuration file.

Host anywhere

MkDocs builds completely static HTML sites that you can host on GitHub pages, Amazon S3, or anywhere else you choose.

Great themes available

There's a stack of good looking themes available for MkDocs. Choose between the built in themes: [mkdocs](#) and [readthedocs](#), select one of the 3rd party themes in the [MkDocs wiki](#), or [build your own](#).

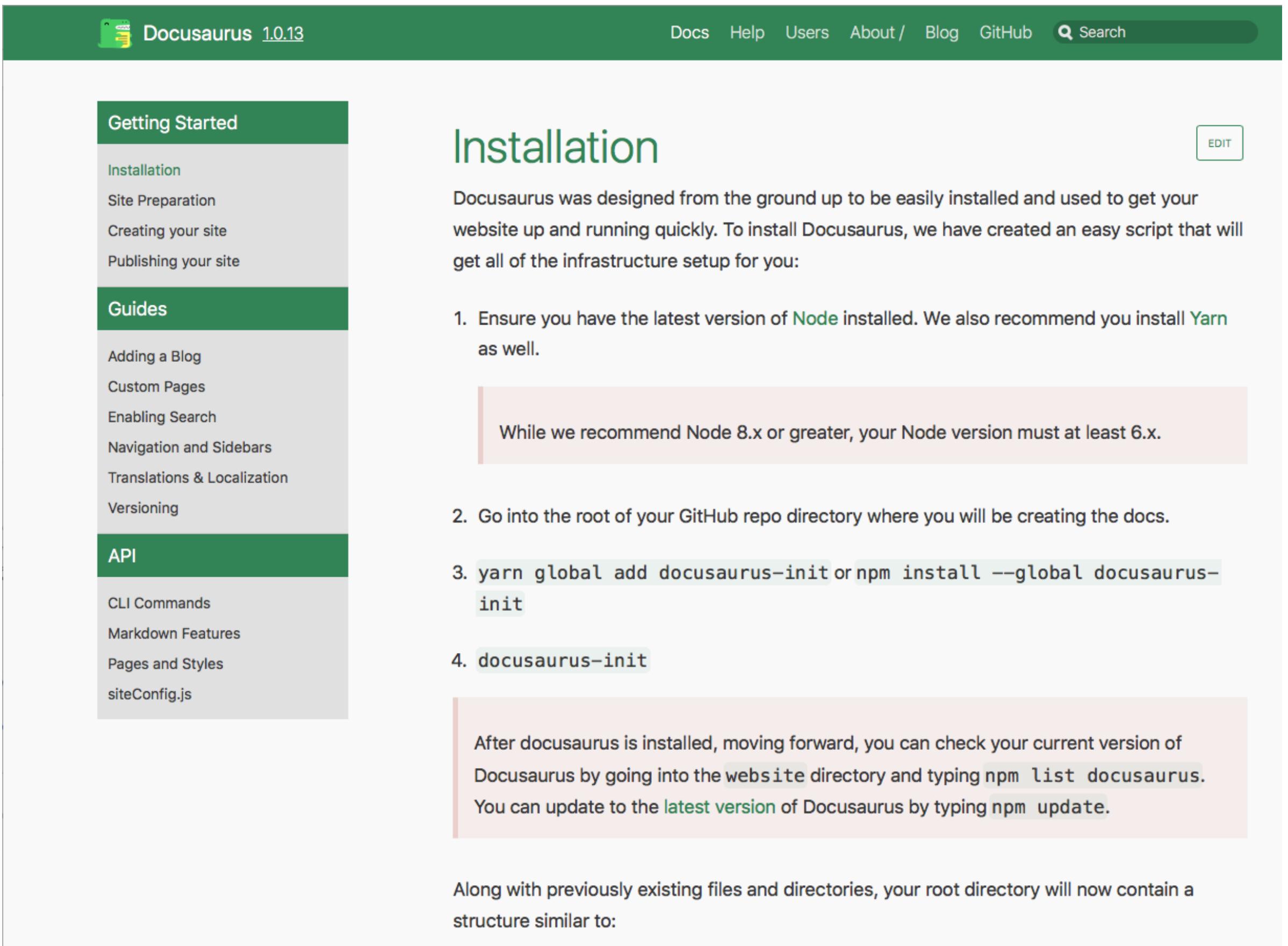
Preview your site as you work

The built-in dev-server allows you to preview your documentation as you're writing it. It will even auto-reload and refresh your browser whenever you save your changes.

Easy to customize

Get your project documentation looking just the way you want it by customizing the theme.

Installation



The screenshot shows the Docusaurus documentation site. At the top, there's a navigation bar with links for Docs, Help, Users, About / Blog, GitHub, and a search bar. The main content area features a sidebar with sections for Getting Started, Guides, and API. The main content area includes a section for Installation with instructions and a note about Node.js requirements, followed by steps for installation and a note about updating.

Docusaurus 1.0.13

Docs Help Users About / Blog GitHub Search

Getting Started

- Installation
- Site Preparation
- Creating your site
- Publishing your site

Guides

- Adding a Blog
- Custom Pages
- Enabling Search
- Navigation and Sidebars
- Translations & Localization
- Versioning

API

- CLI Commands
- Markdown Features
- Pages and Styles
- siteConfig.js

Installation

While we recommend Node 8.x or greater, your Node version must at least 6.x.

1. Ensure you have the latest version of [Node](#) installed. We also recommend you install [Yarn](#) as well.
2. Go into the root of your GitHub repo directory where you will be creating the docs.
3. `yarn global add docusaurus-init` or `npm install --global docusaurus-init`
4. `docusaurus-init`

After docusaurus is installed, moving forward, you can check your current version of Docusaurus by going into the `website` directory and typing `npm list docusaurus`. You can update to the latest version of Docusaurus by typing `npm update`.

Along with previously existing files and directories, your root directory will now contain a structure similar to:

Reproducible Software

- Version Control
- Ease of use
 - Readable
 - Modular
 - Documentation
- **Rigor**
 - Testing
 - CI
- Packaging code
 - Cross platform
 - Containerization

What can happen when you don't test?

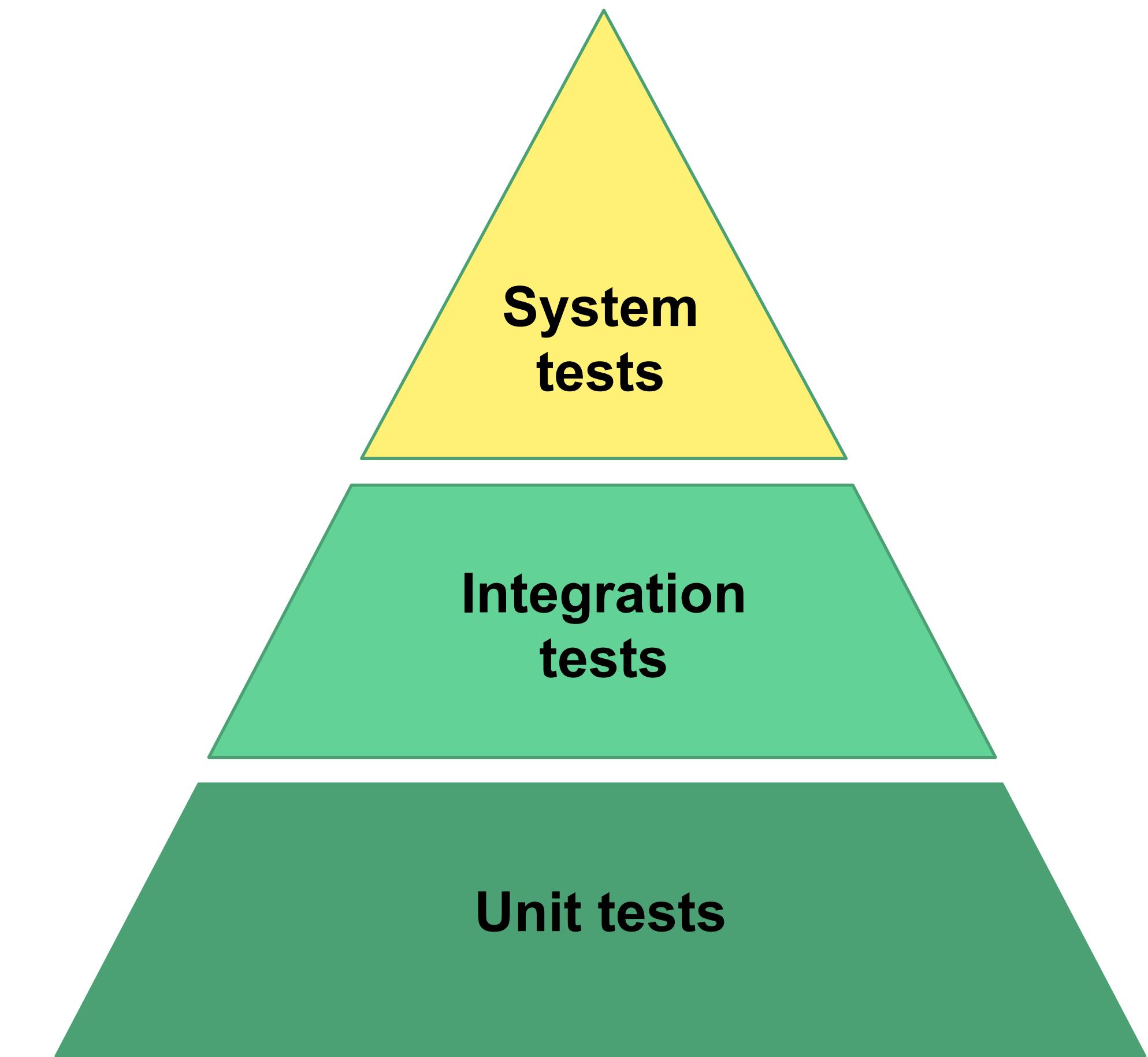


- Navigation software for Ariane-5 had been untested
- Software converted 64-bit floating-point data into 16-bit unsigned integer values
- Arithmetic overflow



Testing

- Prevent bugs in software before they happen
- Make sure discovered bugs don't happen again
- Make major updates go more smoothly, i.e. less debugging



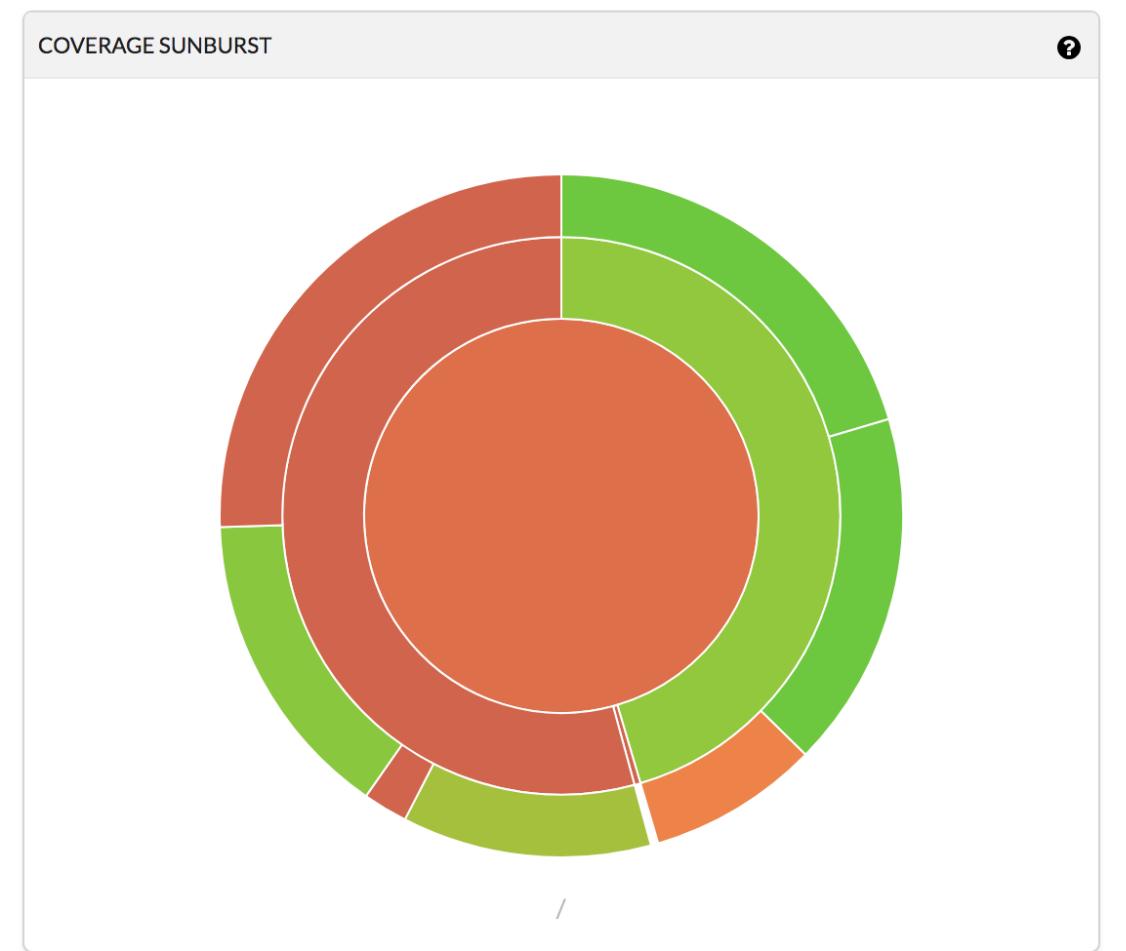
pytest



Testing

Unit tests

- Test each function/component

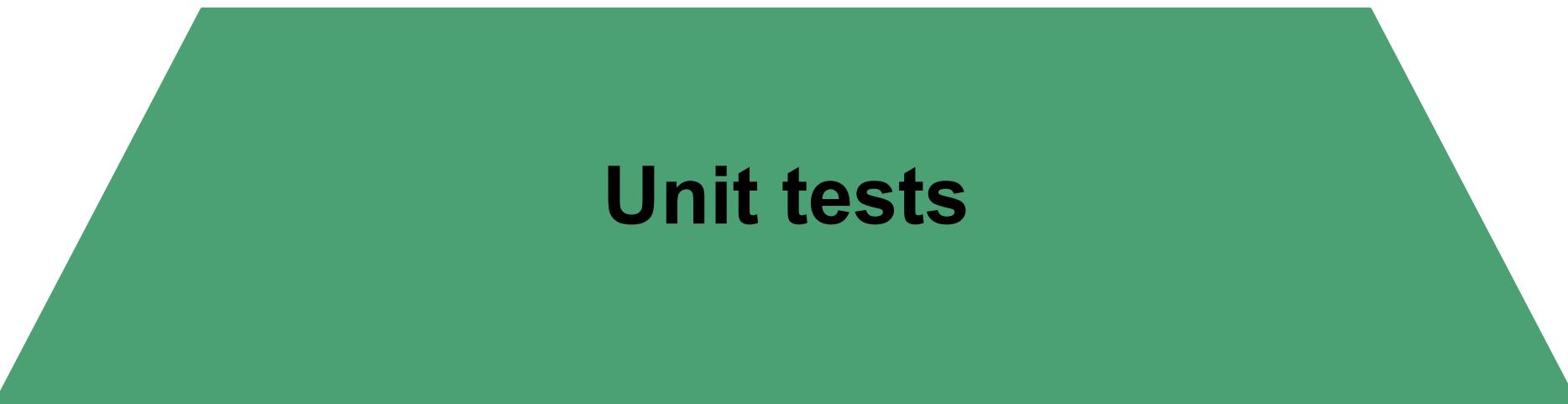


ALL RECENT COMMITS

User	Comment	Date	Issue	Commit Hash	CI Status	Action
aguang	am not sure why we needed adegenet except maybe for outbreaker	5 months ago	issue-8	5a7d9b8	CI Passed	Browse Report
aguang	transmission tree should have been using nmut not onset	5 months ago	issue-8	8e8af85	CI Passed	Browse Report
aguang	updated transmission to use nmut instead of onset	5 months ago	issue-8	1ba01be	CI Passed	Browse Report
aguang	should work now; simpphy_lnx64 with permissions	5 months ago	issue-8	289cf41	CI Passed	Browse Report
aguang	added test for transmission and viral sequences having same individuals	6 months ago	issue-8	c67c482	CI Passed	Browse Report
aguang	forgot test yaml files	6 months ago	issue-8	be6c786	CI Passed	Browse Report

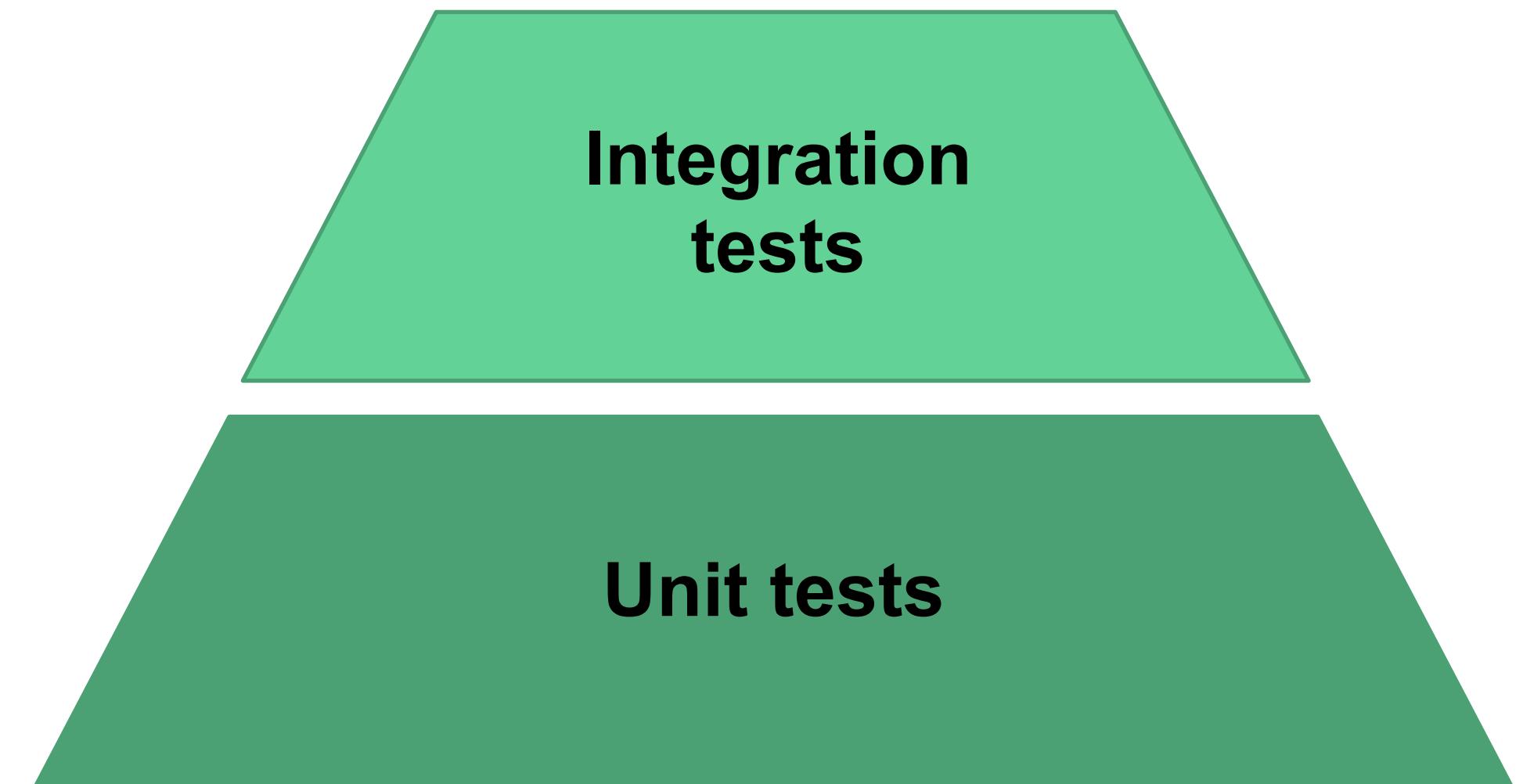
[View all recent commits on issue-8](#)

Files	Lines	Covered	Partially Covered	Uncovered	Coverage
tests	258	247	0	11	95.73%
transmissim	308	165	0	143	53.57%
setup.py	2	0	0	2	0.00%
Project Totals (8 files)	568	412	0	156	72.53%



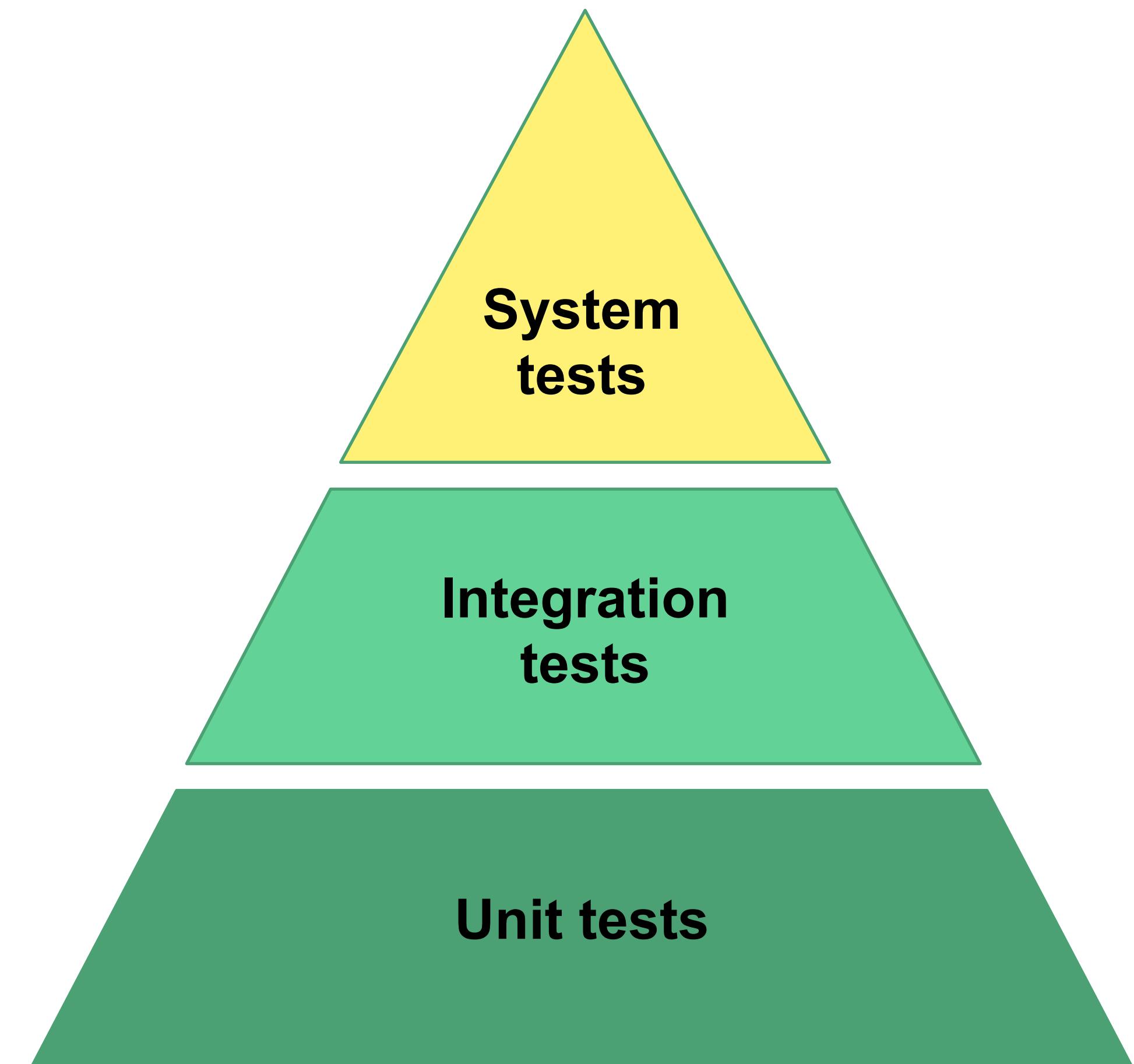
Testing

- Unit tests
 - Test each function/component
- Integration tests
 - Test multiple components together
 - Do they integrate properly with each other?
 - Starts to require test data

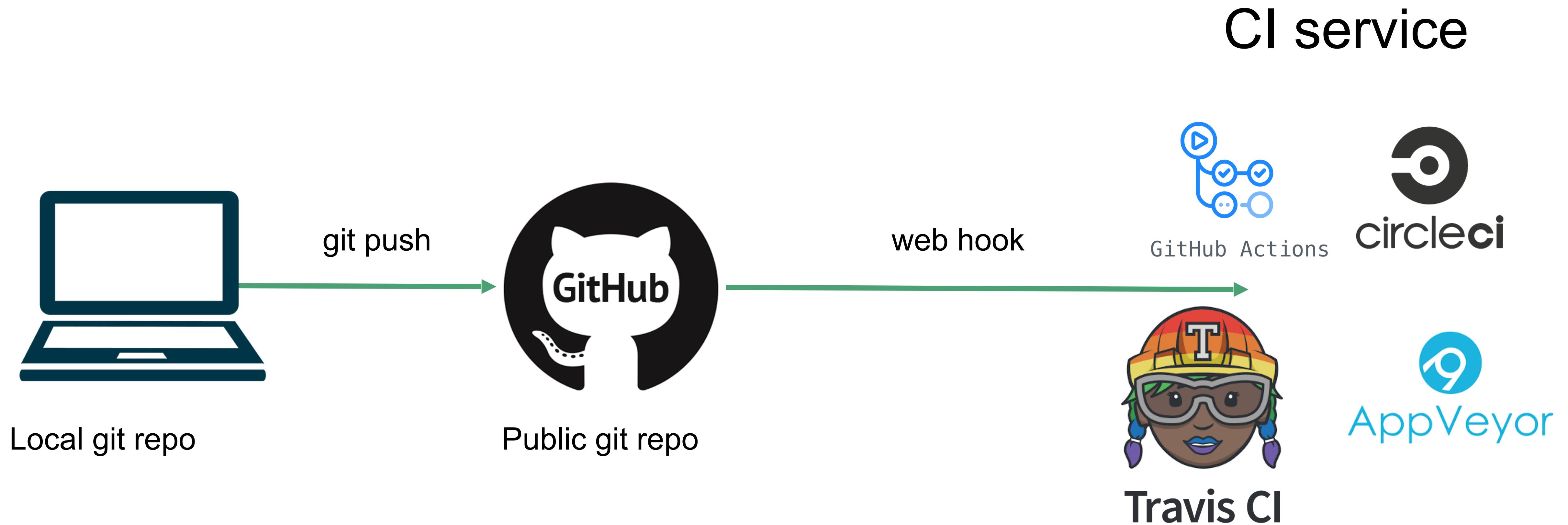


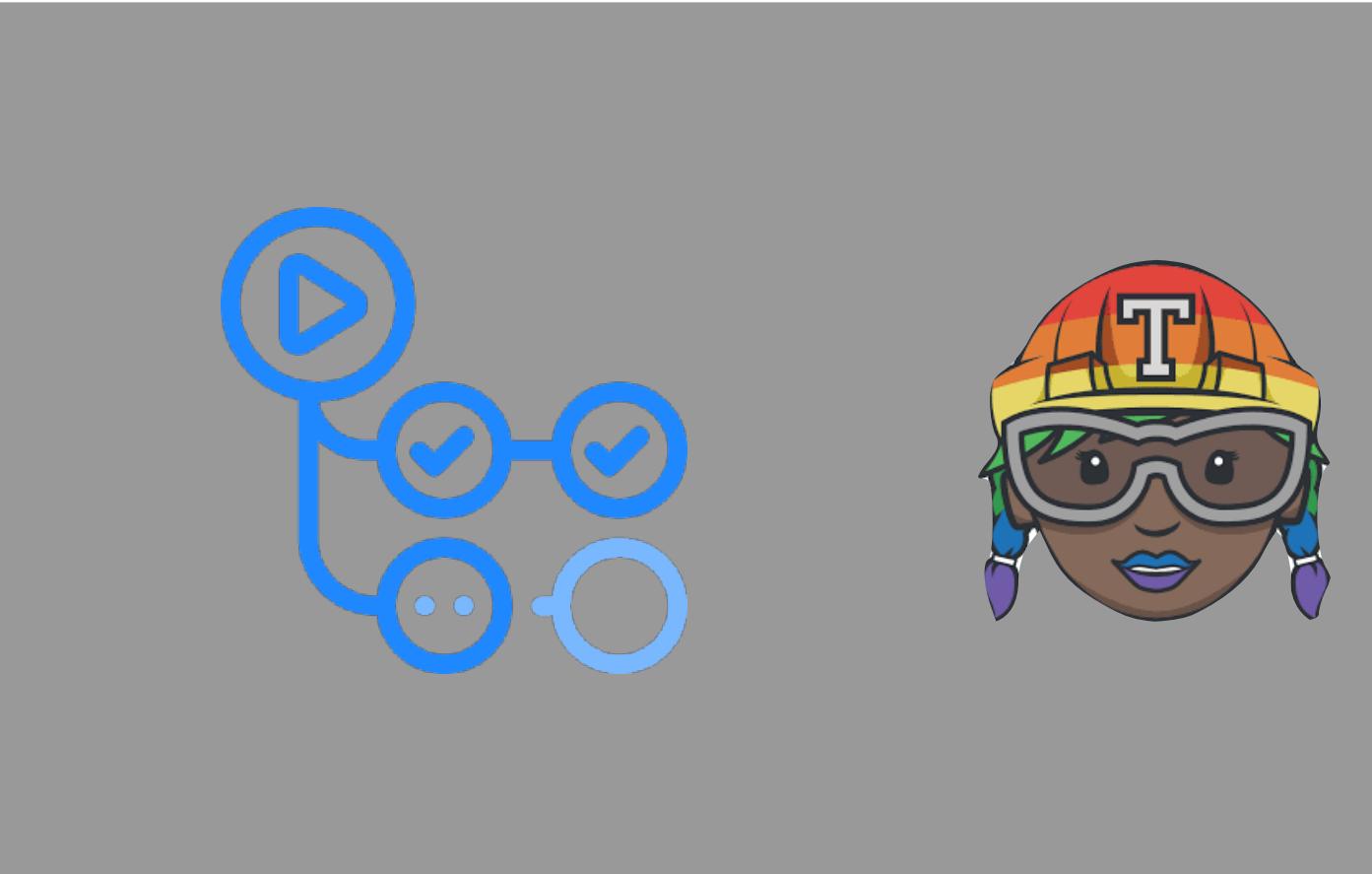
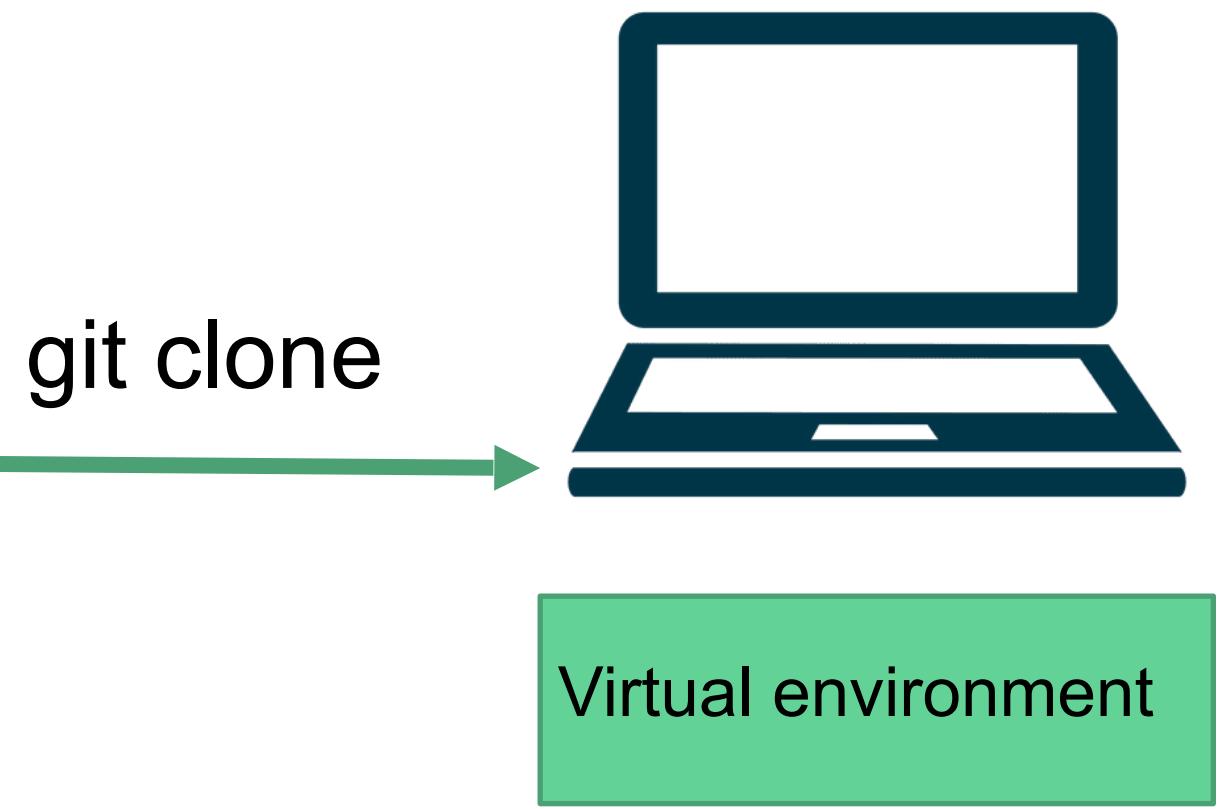
Testing

- Unit tests
 - Test each function/component
- Integration tests
 - Test multiple components together
 - Do they integrate properly with each other?
 - Starts to require test data
- System tests
 - Test entire system
 - Does system work as expected?
- At all levels, should also test that errors are handled correctly!

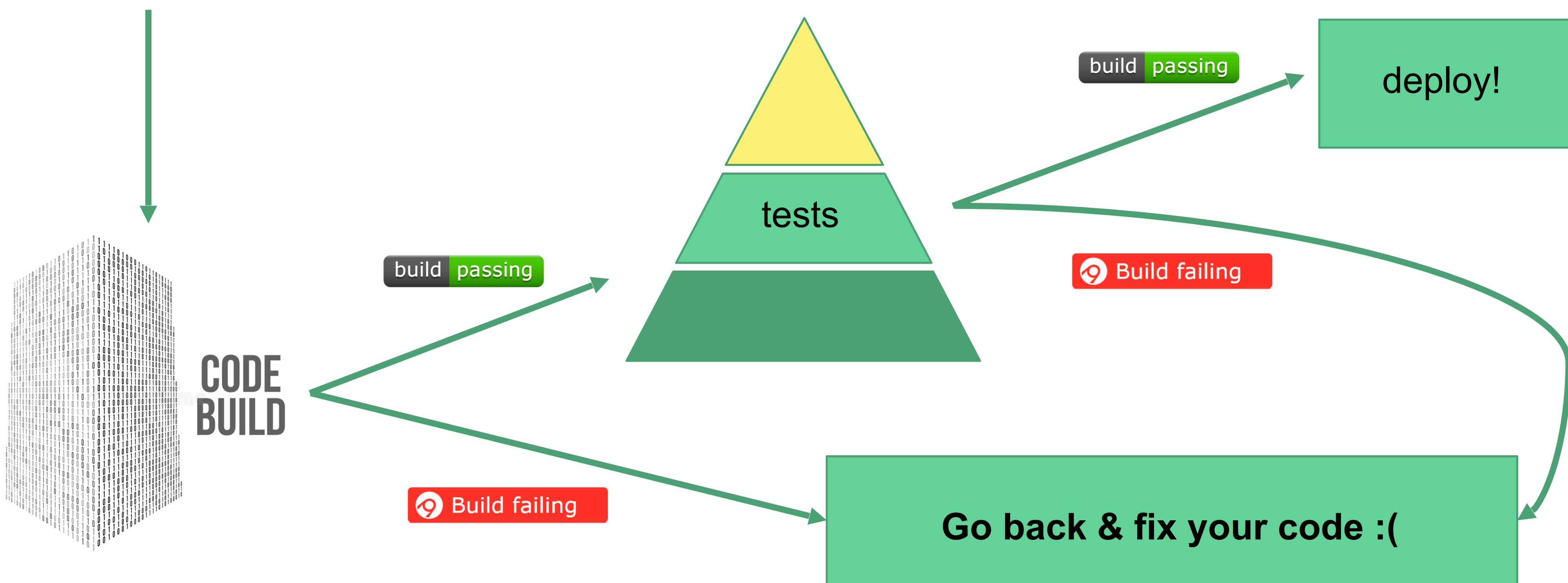


Continuous Integration





CI services automate all of this



Best Practices

- Version Control
- Ease of use
 - Readable
 - Modular
 - Documentation
- Rigor
 - Testing
 - CI
- **Packaging code**
 - Cross platform
 - Containerization

Packages

- Fundamental units of reproducible code
- Includes (ideally):
 - Functions
 - Documentation on how to use them
 - Sample data
 - Unit tests
 - Stable releases
- Often packages on repositories go through a code review that enforces stylistic standards and improves quality



How to package code?

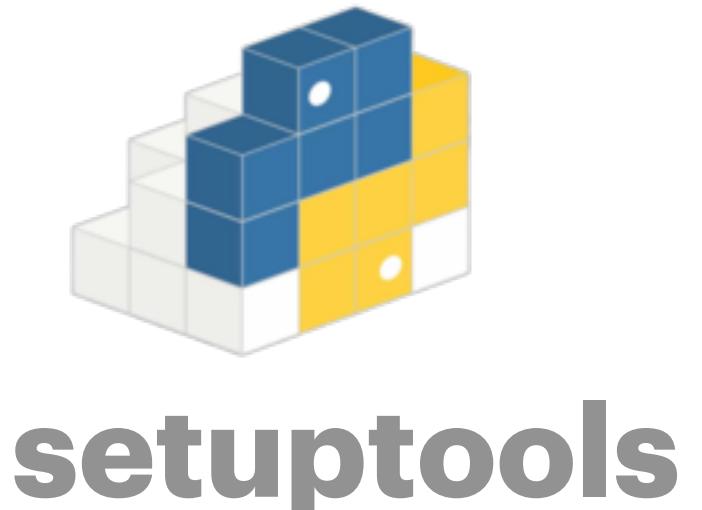
Packaging code

- Local development
 - Virtual environments
 - Dependency management
 - Packaging management

pyenv + Poetry



Welcome to Python's HELL





Open-source cross-platform
package & environment manager

Analysis 1

Analysis 2

Analysis 1

conda env 1

conda env 2

conda env 1

Computer 1

Computer 2

Data science's favorite



- Very popular for data science
- Requires packages to be available in CONDA's registry
- Can install PYPI installs (can create conflicts) `pip install pkg`
- Typically people publishes first to PYPI and then to CONDA
- It's more than a package manager, it can link to different low-level libraries (libc)
- Can be slow (recent improvements and c++ alternatives)

General package manager





pyenv

Installation

List versions of Python available to install

Install most recent stable version

Activate global environment

Listing available versions to the system

Configure a local environment (.python-version)

pyenv manages your
python versions

`brew install pyenv`

`pyenv install --list`

`pyenv install 3.8.1`

`pyenv global 3.8.1`

`pyenv versions`

`pyenv local 3.8.0`



Installation

pipX allows you to install Python CLI utilities in their own environment

```
python -m pip install pipx
```

Poetry

Installation

Create a new project

Use virtual environments inside the project's root directory

Making sure that Poetry is using Pyenvs python

Poetry handles dependency and virtual-environment-management in a way that's very intuitive.

```
pipx install poetry
```

```
poetry new poetry-demo
```

```
poetry config virtualenvs.in-project true
```

```
poetry env use $(pyenv which python)
```

Poetry

```
poetry new poetry-demo
```

```
$ poetry add pendulum
```

```
poetry-demo
├── pyproject.toml
├── README.rst
└── poetry_demo
    └── __init__.py
tests
└── __init__.py
    └── test_poetry_demo.py
```

```
[tool.poetry]
name = "poetry-demo"
version = "0.1.0"
description = ""
authors = ["Sébastien Eustace <sebastien@eustace.io>"]

[tool.poetry.dependencies]
python = "*"

[tool.poetry.dev-dependencies]
pytest = "^3.4"
```

```
poetry install
```

If there is no `poetry.lock` file present, Poetry resolves all dependencies listed in your `pyproject.toml` file and downloads the latest version of their files.



Poetry

```
$ poetry build  
  
Building poetry (1.0.0)  
- Building sdist  
- Built poetry-1.0.0.tar.gz  
  
- Building wheel  
- Built poetry-1.0.0-py2.py3-none-any.whl
```

```
$ poetry publish
```

Publishing poetry (1.0.0) to PyPI

- Uploading poetry-1.0.0.tar.gz 100%
- Uploading poetry-1.0.0-py2.py3-none-any.whl

Detour: GitHub Actions Example

<https://github.com/python-poetry/poetry>

Reproducible Software

- Source code
- Third-Party libraries
- Runtime environment

Source Code

- Captured in version control
- Releases tagged
- Immutable tags

Third-Party Libraries

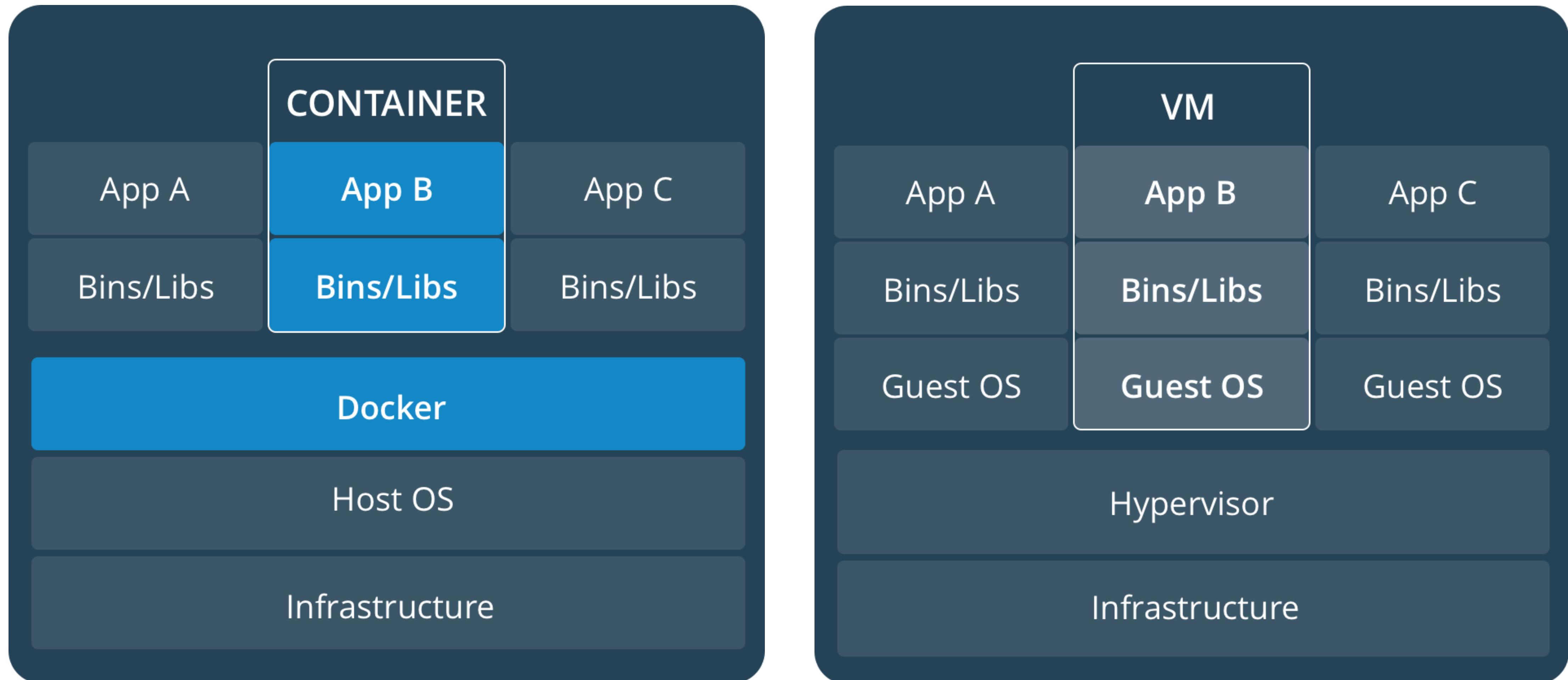
- Captured in a package manifest file
- Libraries pegged to specific version
- Peg libraries to one version, not a range
- Select libraries that are maintained and have a corporate sponsor

Runtime environment

- Capture system software requirements
- Freeze runtime environment
- Package runtime environment for consumption



Docker vs. VM



source: docs.docker.com

Docker - Creating an Image

There are several ways to create a Docker image:

- dockerfile (write code yourself)
- repo2docker (create an image from a github repo)

Once created, images can be shared, either using the central ‘DockerHub’ repository, GitHub or by exporting and sharing a tarball

Dockerfile

-FROM python

+FROM python:3.8

```
RUN apt-get update && apt-
get install -y \
    libxml2-dev \
```

-FROM python:3.8

+FROM python:3.8.0-buster

```
RUN apt-get update && apt-
get install -y \
    libxml2-dev \
```

Always specify full version
for images

Dockerfile

```
FROM python:3.8.0-buster
```

```
RUN apt-get update && apt-get  
install -y \  
- libxml2-dev \  
- libssqlite3-dev  
+ libxml2-dev=2.9.4+dfsg1-7+b3  
\  
+ libssqlite3-dev=3.27.2-3
```

```
COPY requirements.txt .  
RUN pip install -r  
requirements.txt
```

Always specify full version for images

Always specify versions of system software

Dockerfile

```
WORKDIR /app  
COPY . ./  
  
+ENV RANDOM_SEED=1024  
ENTRYPOINT ["python"]  
-CMD ["my-script.py"]  
+CMD ["my-script.py", "special",  
"arguments", "here"]
```

Always specify full version
for images

Always specify versions of
system software

Capture arguments and env
variables

Gotchas

- Understand the OS your image uses
- Different operating systems have different packages
- Alpine linux uses a different libc (glibc vs musl)

Building Containers

Build from Dockerfile in current directory

Build and give a name

Build, name and tag

```
docker build .
```

```
docker build -t science
```

```
docker build -t science:1.0.7
```

Reproducible Data

- Make it available!
- Keep track of transformations, data-provenance



Git Large File Storage

Installation

Set up git-lfs for your user

Indicate which files git-lfs should track

Add .gitattributes file

Add, commit and push as usual

GIT-LFS or GIT ANNEX allow you to version control your data

```
brew install git-lfs
```

```
git lfs install
```

```
git lfs track "*.psd"
```

```
git add .gitattributes
```

```
git add file.psd
```

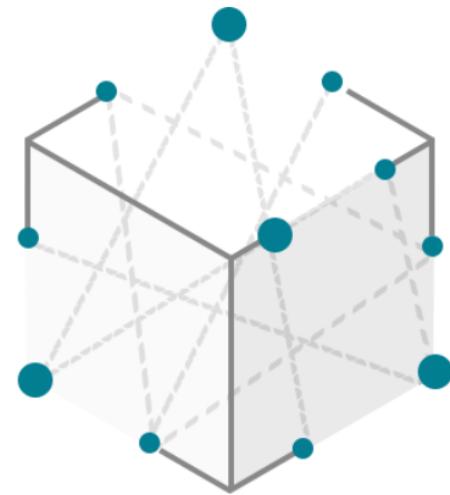
```
git commit -m "Add design file"
```

```
git push origin master
```

kaggle



Leverage domain-specific data repositories



OpenNEURO



GIN

Modern Research Data Management for Neuroscience

...distributed version control, flavoured for science

Snakemake

nextflow

Reproducible Workflows



Snakemake

```
# Count words.
rule count_words:
    input: 'books/isles.txt'
    output: 'isles.dat'
    shell: 'python wordcount.py
books/isles.txt isles.dat'
```

- Keep track of all the steps performed in an analysis
- Bash scripts are not an efficient way of storing a workflows

Reproducible Data + Workflow + Parameters



Data
Lad

The word "Data" is in a black, lowercase, sans-serif font, while "Lad" is in a bright orange, lowercase, sans-serif font. The "L" in "Lad" has a small orange arrow pointing upwards at its top-right corner.



OPEN-SOURCE
MACHINE LEARNING
VERSION CONTROL

Connect code and data by commands

Make changes and reproduce

Share code

Share data and ML models

Track data

```
git add train.py  
dvc add images.zip
```

```
dvc run -d images.zip -o images/ unzip  
-q images.zip  
dvc run -d images/ -d train.py -o  
model.p python train.py
```

```
vi train.py  
dvc repro model.p.dvc
```

```
git add .  
git commit -m 'The baseline model'  
git push
```

```
dvc remote add myremote -d s3://  
mybucket/image_cnn  
dvc push
```

Questions

Phew! That was a lot!