# Introduction to Workflow Management Solutions for Public Health Bioinformatics

Model for Distributed Public Health Bioinformatics:
Week 3 – Connecting WDL Workflows with Terra.Bio

Tuesday May 11th, 2022
Kevin G. Libuit, MS | Theiagen Genomics

# Training Workshop Overview

**Communication and Support**

- Slack workspaces:
  - **Terra-US-PHL; #wdl-writing**
  - **StaPH-B; #workflow-management, #cromwell_noobz**
- Email: **support@terrapublichealth.zendesk.com**
- Weekly Office Hours:
  - **Mountain Region - Friday 9-10AM (PDT)**
  - **North East Region - Friday 10-11AM (PDT)**

# Training Workshop Overview

**Main Course Objective**

Learn how to use **workflow management systems** to **develop accessible, interoperable, and reproducible** public health bioinformatics solutions

# Last Week's Content: WDL Tasks & Workflows

## Major Takeaways (WDL Task Files):

- **Input Section**
  - **Obligate Inputs:** Required for the task to run successfully
  - **Optional Inputs:** Not required for the task to run successfully
  - **Declared Inputs:** Default values of obligate or optional inputs; can be overridden when task is called in workflow
- **Command Section**
  - Two Options Available to Define a Command Element: <<< >>> or { }
- **Output Section**
  - **Obligate Outputs:** Required for the task to run successfully
  - **Optional Outputs:** Not required for the task to run successfully
- **Runtime Section**
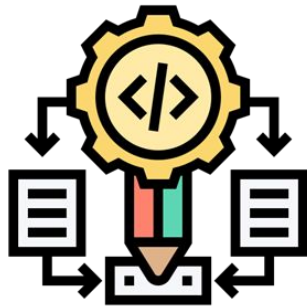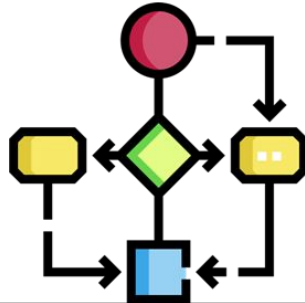  - Recognized runtime attributes depend on **both workflow engine & compute backend**

# Last Week's Content: WDL Tasks & Workflows

## Major Takeaways (WDL Workflow Files):

- **Input Section**
  - **Obligate Inputs:** Required for the workflow to run successfully
  - **Optional Inputs:** Not required for the workflow to run successfully
  - **Declared Inputs:** Default values of obligate or optional inputs; can be overridden when workflow is run
- **Call Section**
  - **Task Inputs:** Declared for each task called – even if no input values are required
  - **Task Aliases:** Helpful when a single task is called multiple times
- **Output Section**
  - **Obligate Outputs:** Required for the task to run successfully
  - **Optional Outputs:** Not required for the task to run successfully

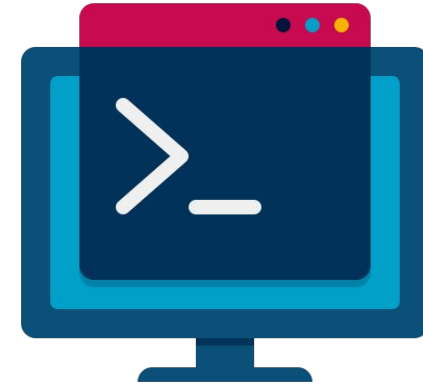# Executing WDL Workflows



**WDL Source Code**

WDL Task
Individual step performing a single and specific bioinformatics job

WDL Workflow
Software that strings together multiple analytical modules
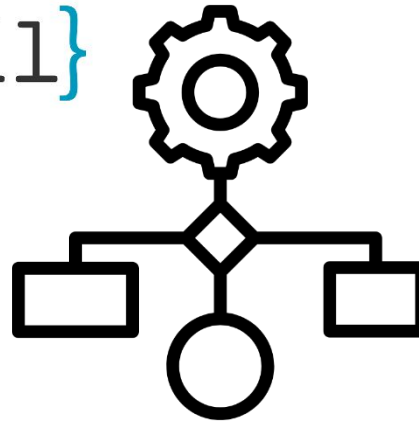
Command-Line Interface (CLI)

Terra

Graphic-User Interface (GUI)

# Introduction to the Terra Platform

Terra is a **bioinformatics web application** that connects users to bioinformatics workflows (WDL) & dynamic cloud computing resources (GCP) through a clean and intuitive user interface

WDL workflows can be accessed on Terra and run on GCP resources

**Terra compatibility expanding:** Will support Nextflow workflows and Azure backend

**Bioinformatics workflows**
Specialized software written to analyze biological data

**Cloud Computing Resources**
Network of remote compute resources hosted on the internet

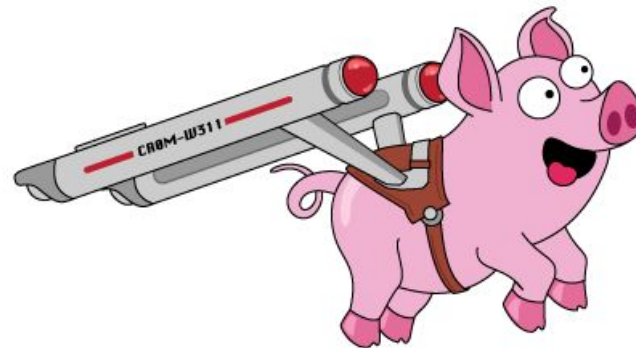# Introduction to the Terra Platform

Terra users can access WDL workflows hosted on Dockstore
- Workflow repository that enables "*researchers and developers to **share and reuse** analytical **workflows and tools** in a way that makes them **machine readable and runnable** in a **variety of environments**"*

From Terra, workflows are executed using the **Cromwell engine**
- Execution on GCP backend

# Running Scan-N-Trim Workflow



Command-Line Interface (CLI)

Graphic-User Interface (GUI)

# Connecting WDL Workflows to Terra.Bio



**Publicly Accessible
Github Repository**

Containerized WDL Workflows

**Linked Dockstore Account**

Can continuously configure on GitHub
repository after initial setup

# Connecting WDL Workflows to Terra.Bio



**WDL Source Code**

**Hosting a Terra-Accessible WDL Workflow on Dockstore**
1. Host WDL source code on publicly accessible GitHub repository
2. Create a Dockstore account
3. Sync GitHub repo with Dockstore account
4. Create a .dockstore.yml file into the top level of GitHub repository
5. Publish workflow from Dockstore

# Terra Platform – Data Tables & Workflows

**Terra Data Table**

How data is organized within the Terra platform

- Data within a Terra Data Table can be uploaded by the user or generated by a Terra Workflow

**Terra Workflows**

WDL workflows executable from the Terra platform

- Can take inputs from a Terra Data Table

- Populates outputs to a Terra Data Table

**WORKSPACES** BETA

Terra

Workspaces › theiagen-validations/wm_training_klibuit ›
**Data**

DASHBOARD | DATA

Data tab is selected to view all Terra Data Tables available in this workspace

**TABLES** ⊕

📖 **wm_training_specimen** (1)

**REFERENCE DATA** ⊕

**OTHER DATA**

📖 Workspace Data

Data within the wm_training_specimen Terra Data Table is displayed

⬇ DOWNLOAD ALL ROWS | 📋 COPY PAGE TO CLIPBOARD | 0 rows se

Data within any Terra Data Table can be analyzed with a Terra Workflow

| ☐ ▼ | wm_training_specimen_id ⓘ ⬇ | read1 ⓘ | read2 ⓘ |
|---|---|---|---|
| ☐ | sample_01 | sample01_R1.fastq.gz | sample01_R1.fastq.gz |

Each row within a Terra Data Table is called an *entity* each column represents some *attribute* associated with any given entity

For the purpose of this training, **data must be organized into a Terra Data Table** to be analyzed

# Terra Platform – Data Tables & Workflows

## Terra Data Table

How data is organized within the Terra platform

- Data within a Terra Data Table can be uploaded by the user or generated by a Terra Workflow

## Terra Workflows

Bioinformatics Workflows executable from the Terra platform

- Takes input from a Terra Data Table

- Populates outputs to a Terra Data Table

← Back to list

## Scan-N-Trim

Version: | solutions ▾ |

Source: github.com/theiagen/wm_training/Scan-N-Trim:solutions

Synopsis:

*No documentation provided*

○ Run workflow with inputs defined by file paths

● Run workflow(s) with inputs defined by data table

**Step 1**

Select root entity type: | wm_training_spec... ▾ |

**Step 2**

| SELECT DATA |  No data selected

☑ Use call caching     ☐ Delete intermediate outputs ⓘ     ☐ Use reference disks ⓘ     ☐ Retry with more memory ⓘ

SCRIPT     • •     **INPUTS**     • •     OUTPUTS     • •     RUN ANALYSIS

← Back to list

# ⦂ Scan-N-Trim

Version: | solutions ⌄ |  ⎫ Can select GitHub branch or
                          version tag

Source: github.com/theiagen/wm_training/Scan-N-Trim:solutions

Synopsis:

*No documentation provided*

○ Run workflow with inputs defined by file paths

● Run workflow(s) with inputs defined by data table

**Selecting the Terra Data Table** that
contains the input data to be analyzed

┌ Step I

Select root entity type: | wm_training_spec... ⌄ |     **SELECT DATA**    No data selected

☑ Use call caching    ☐ Delete intermediate outputs ⓘ    ☐ Us        INPUT & OUTPUT forms to    more memory ⓘ
                                                                 **define the input data to analyze
                                                                  and what outputs to generate**

SCRIPT    • •    INPUTS    • •    OUTPUTS    • •    RUN ANALYSIS

SCRIPT    • •    INPUTS    • •    ANALYSIS

INPUTS form selected

Hide optional inputs

| Task name ↓ | Variable | Type | Attribute |
|---|---|---|---|
| scan_n_trim_workflow | read1 | File | this.read1 |
| scan_n_trim_workflow | read2 | File | this.read2 |

Inputs defined using "*this.{attribute}*" notation
*this* = selected root entity
*{attribute}* = attribute within proceeding item
*this.reads* =  reads attribute within the root entity

SCRIPT ●● INPUTS ●● OUTPUTS ●●

Output files will be saved to

📁 Files / *submission unique ID* / scan_n_trim_workflow / *workflow unique ID*

References to outputs will be written to

📄 Tables / wm_training_specimen

Fill in the attributes below to add or update columns in your data table

Select "Use defaults" to ensure all workflow outputs are populated to your Terra table

| Task name ↓ | Variable | Type | Attribute  \|  Use defaults |
|---|---|---|---|
| scan_n_trim_workflow | read1_clean_total_reads | Int | this.read1_clean_total_reads |
| scan_n_trim_workflow | read1_raw_total_reads | Int | ds |
| scan_n_trim_workflow | read2_clean_total_reads | Int | this.read2_clean_total_reads |
| scan_n_trim_workflow | read2_raw_total_reads | Int | this.read2_raw_total_reads |

Will create a **column** to populate every output defined in the **output section** of the **WDL workflow**

# Connecting WDL Workflows with Terra.Bio

## Major Takeaways:

- Terra is a **bioinformatics web application** that connects users to bioinformatics workflows (WDL) & dynamic cloud computing resources (GCP) through a **clean and intuitive user interface**
  - Terra compatibility expanding to support **Nextflow workflows** and **Azure backend**
- Terra runs **WDL workflows** using the **Cromwell engine** and executes on **GCP resources**
  - Terra account setup requires GCP account
- WDL workflows must be **hosted on Dockstore** to enable Terra access
  - Requires synchronizing a public GitHub repository with Dockstore account

# Lecture Exercise: Running the Scan-N-Trim WDL Workflow From Terra

## [Scan_N_Trim](#):
- Navigate to the Terra workspace provided to you for this training
- From the Workflows tab, select the Scan-N-Trim workflow
- Run this workflow with inputs defined by data table
  - Select the wm_training_specimen as the root entity
  - Define the required inputs and use defaults for the outputs
  - Run this workflow on sample_01

Complete this exercise during our 20m session break;
we will begin again at 12:08PM (PT)

# Lecture Exercise: Linking Your WDL Workflow to Terra

## For Trainees that Have Shared a GitHub Username:
- Create a dev branch on the wm_training repository
- Merge your work onto this branch and push all commits to GitHub
- Write a .dockstore.yml file to host your Week 2 solution to Dockstore and push commit
  - Workflow name must be set as `Scan-N-Trim`
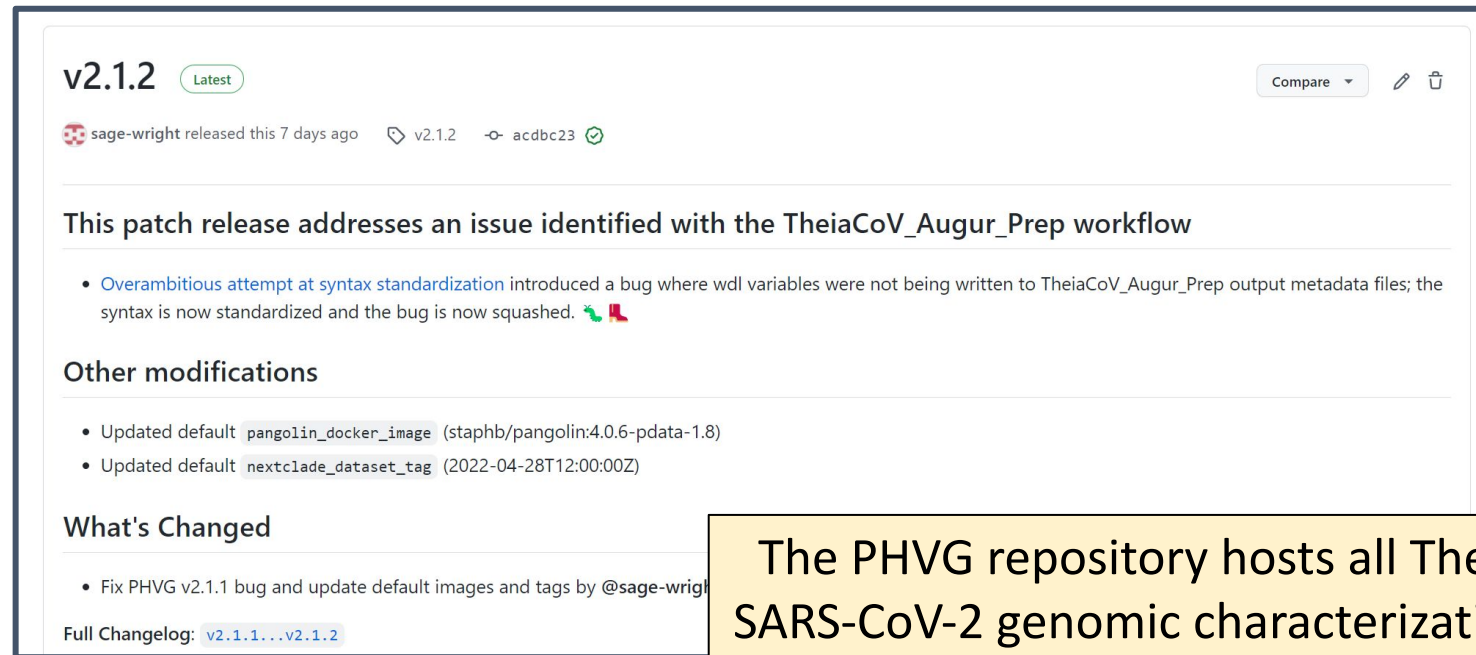- Run this workflow from your Terra workspace

## For Trainees that Have Not Shared a GitHub Username:
- Modify your Week 2 solution to capture the cleaned read files
- Add a task to your Week 2 solution to:
  - Perform Shovill assembly using these cleaned read files
  - Assess the quality of this assembly using QUAST

# Terra Workflow Versions

**Linked to GitHub Version Releases and Branches**
- **Version Release**: Static iteration of a GitHub repository
  - Ideally a validated release of a code base that will not change



The PHVG repository hosts all Theiagen workflows for SARS-CoV-2 genomic characterization, epidemiology and submission preparation

# Terra Workflow Versions

**Linked to GitHub Version Releases and Branches**
- **Version Release**: Static iteration of a GitHub repository
    - Ideally a validated release of a code base that will not change



We recommend that most laboratories **utilize a versioned release** when running workflows on Terra
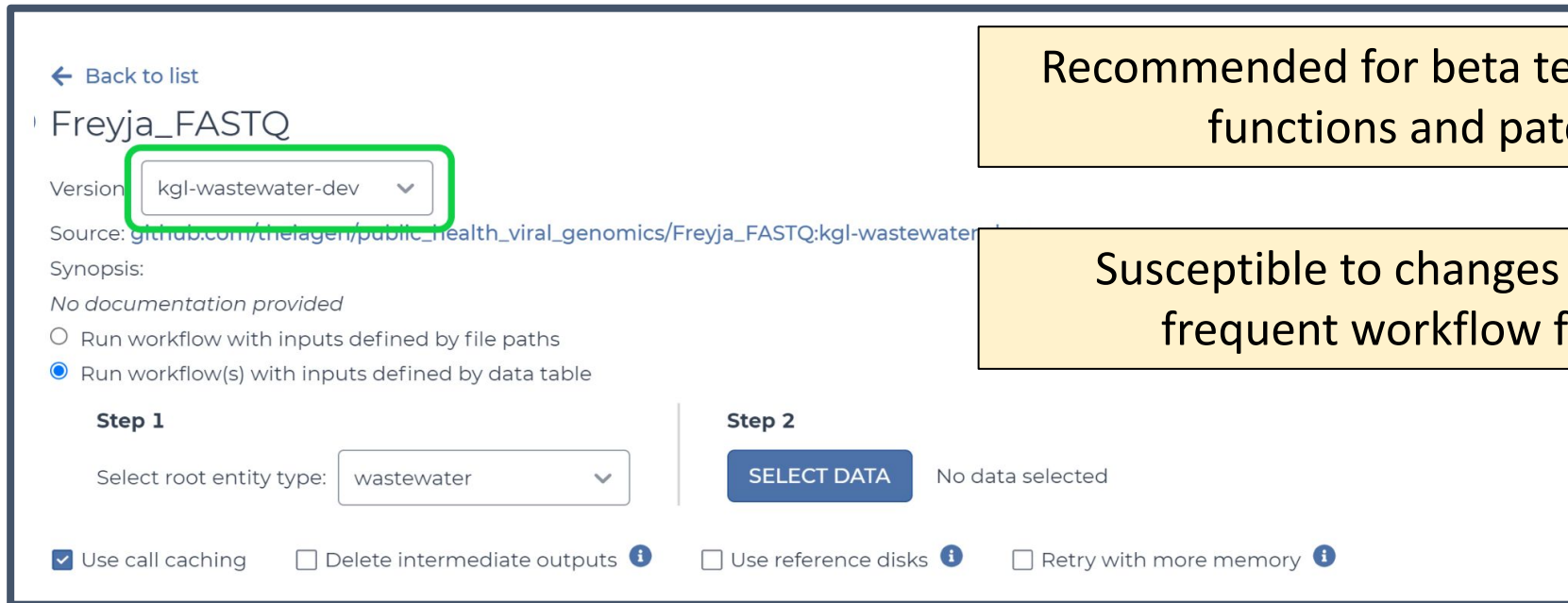
# Terra Workflow Versions

**Linked to GitHub Version Releases and Branches**

- **Branches**: Dynamic iteration of a GitHub repository
  - Unvalidated and under active development; susceptible to more regular code changes for optimization and testing
  - Utilized to validate patches or new workflow functions before including into a versioned release
  - Theiagen practice:
    - Designate the name of the branch with the developer, new function/patch, and "*dev*", e.g. "*kgl-wastewater-dev*"

# Terra Workflow Versions

## Linked to GitHub Version Releases and Branches

- **Branches**: Dynamic iteration of a GitHub repository



← Back to list

Freyja_FASTQ

Version [ kgl-wastewater-dev ▾ ]

Source: github.com/thelagen/public_health_viral_genomics/Freyja_FASTQ:kgl-wastewater-

Synopsis:
*No documentation provided*

○ Run workflow with inputs defined by file paths
◉ Run workflow(s) with inputs defined by data table

**Step 1**                                      **Step 2**

Select root entity type:  [ wastewater ▾ ]     [ SELECT DATA ]  No data selected

☑ Use call caching    ☐ Delete intermediate outputs ⓘ   ☐ Use reference disks ⓘ   ☐ Retry with more memory ⓘ

Recommended for beta testers of new functions and patches

Susceptible to changes and more frequent workflow failures

# Terra Workflow Versions

**Linked to GitHub Version Releases and Branches**
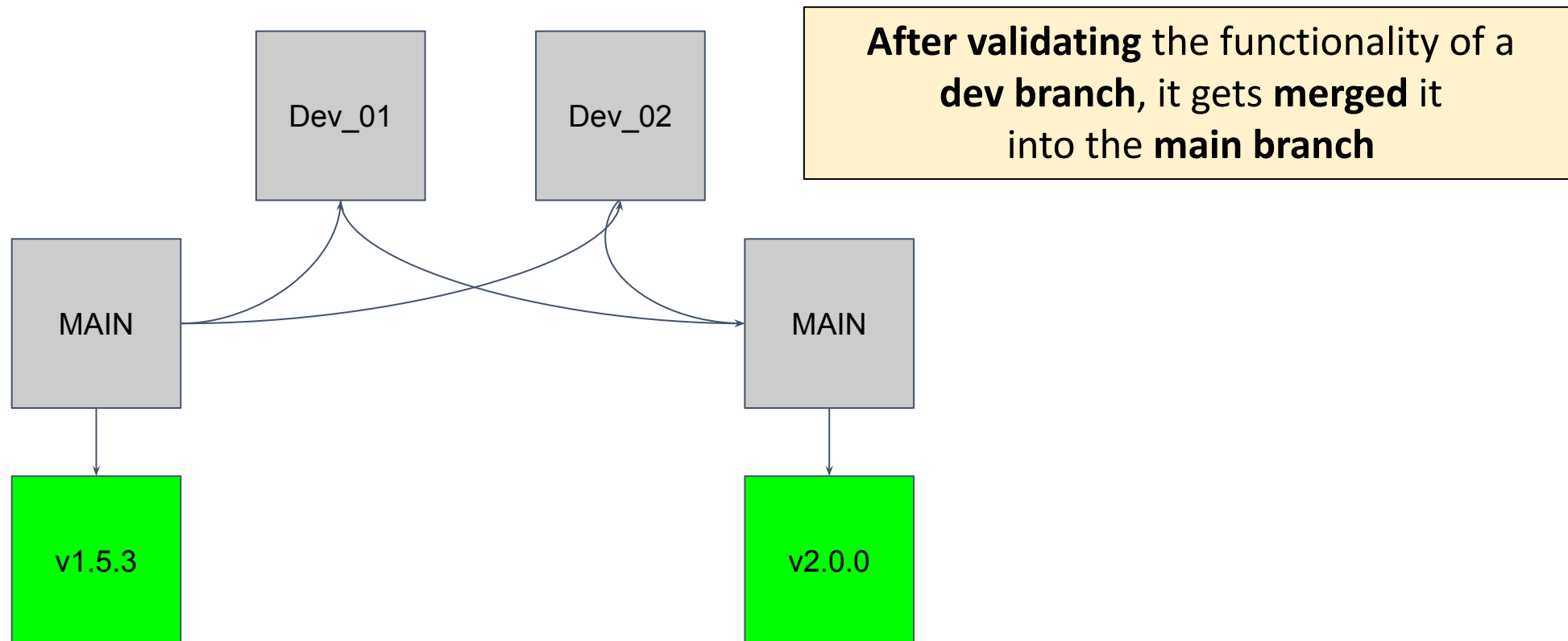- **Branches**: Dynamic iteration of a GitHub repository



The *main* branch is also under active development, but more stable than *dev* branches

# Terra Workflow Versions

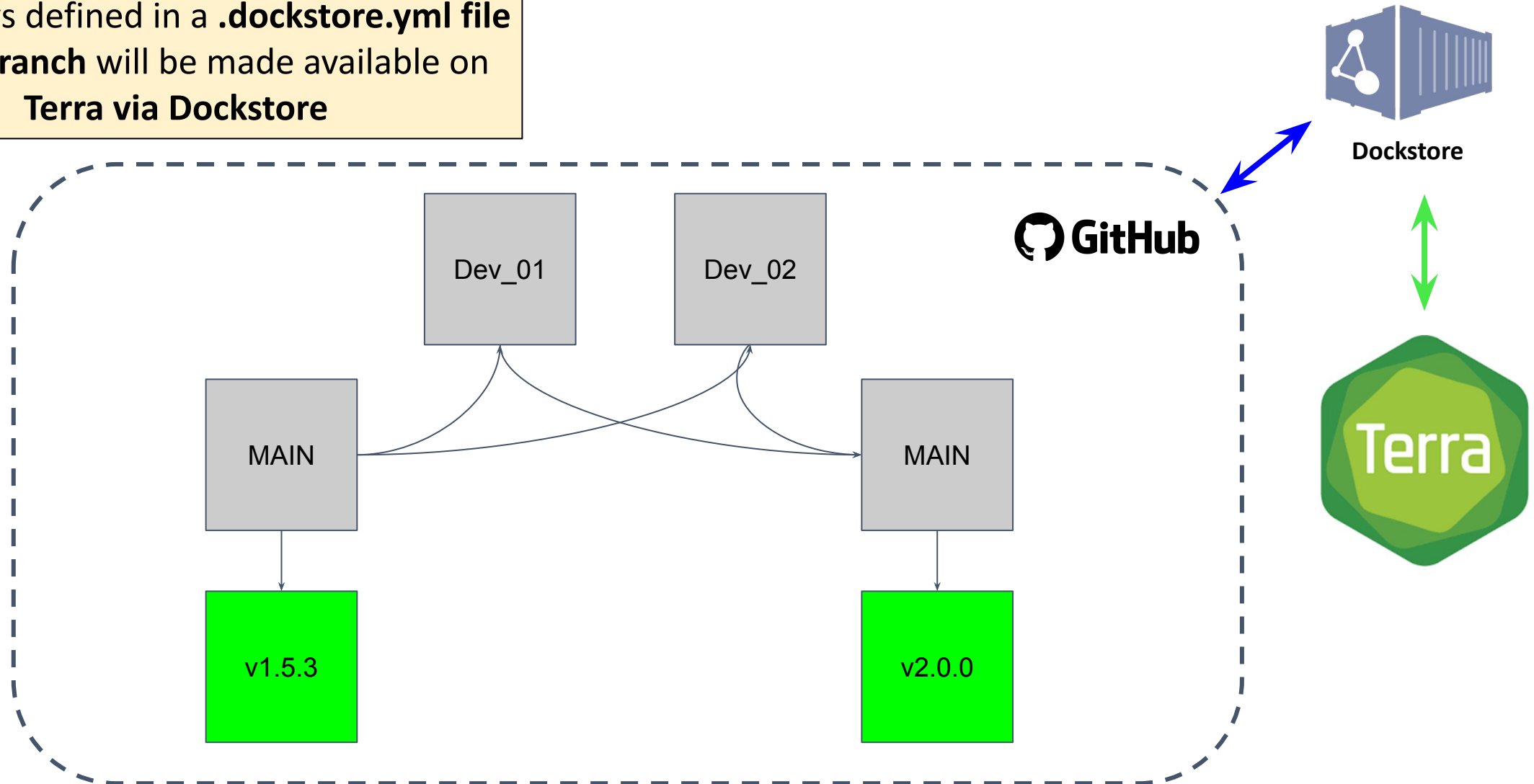**Linked to GitHub Version Releases and Branches**

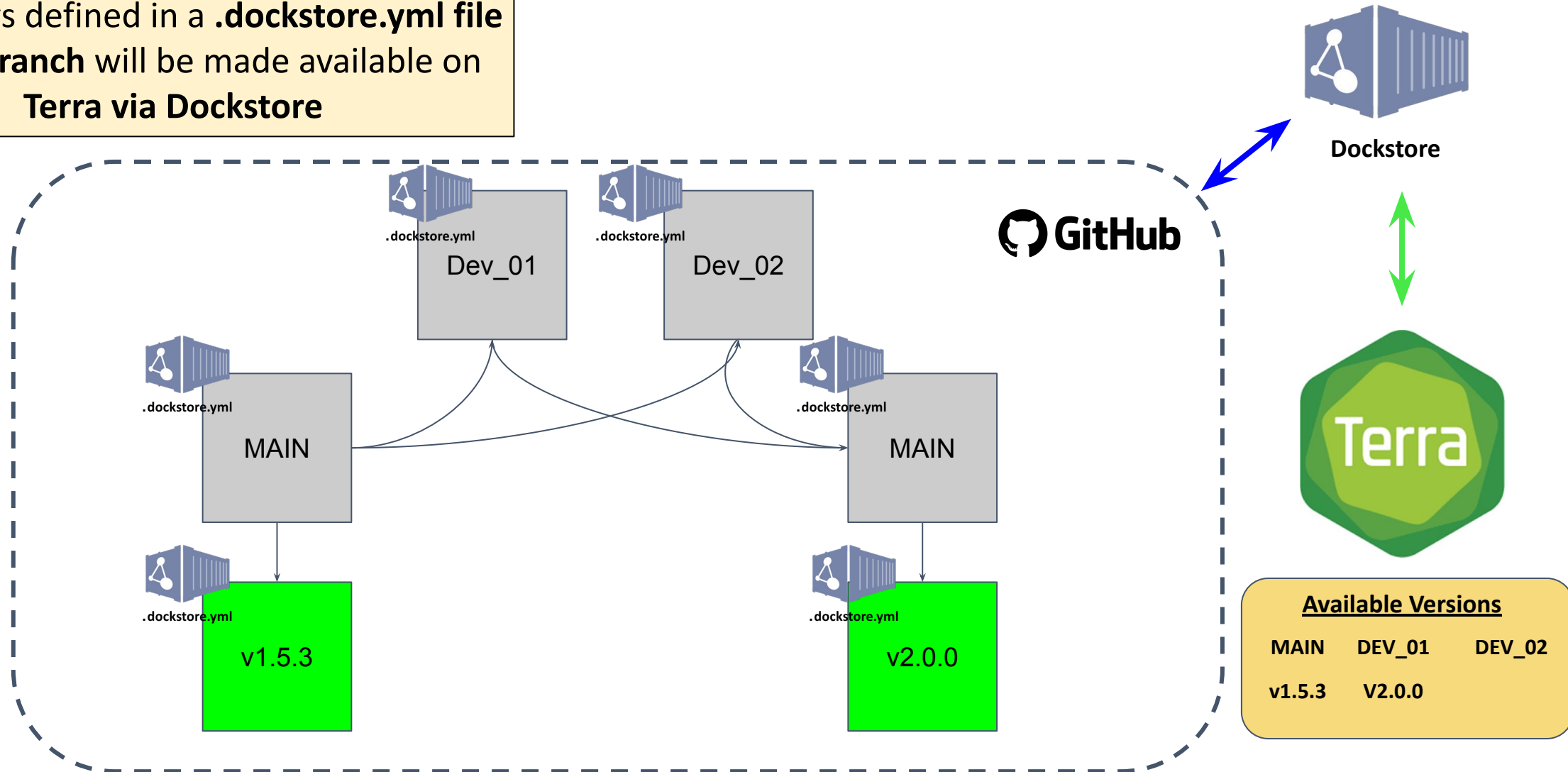- **Branches**: Dynamic iteration of a GitHub repository



> **After validating** the functionality of a **dev branch**, it gets **merged** it into the **main branch**

# Terra Workflow Versions

Workflows defined in a **.dockstore.yml file** on **any branch** will be made available on **Terra via Dockstore**

Dockstore

GitHub

Dev_01

Dev_02

MAIN

MAIN

v1.5.3

v2.0.0

Terra

# Terra Workflow Versions

Workflows defined in a **.dockstore.yml file** on **any branch** will be made available on **Terra via Dockstore**



Dockstore

GitHub

.dockstore.yml
Dev_01

.dockstore.yml
Dev_02

.dockstore.yml
MAIN

.dockstore.yml
MAIN

.dockstore.yml
v1.5.3

.dockstore.yml
v2.0.0

Terra

**Available Versions**

MAIN      DEV_01      DEV_02

v1.5.3      V2.0.0

# Lecture Exercise: Linking Your WDL Workflow to Terra

## For Trainees that Have Shared a GitHub Username:

- Create a dev branch on the wm_training repository
- Merge your work onto this branch and push all commits to GitHub
- Write a .dockstore.yml file to host your Week 2 solution to Dockstore and push commit
  - Workflow name must be set as `Scan-N-Trim`
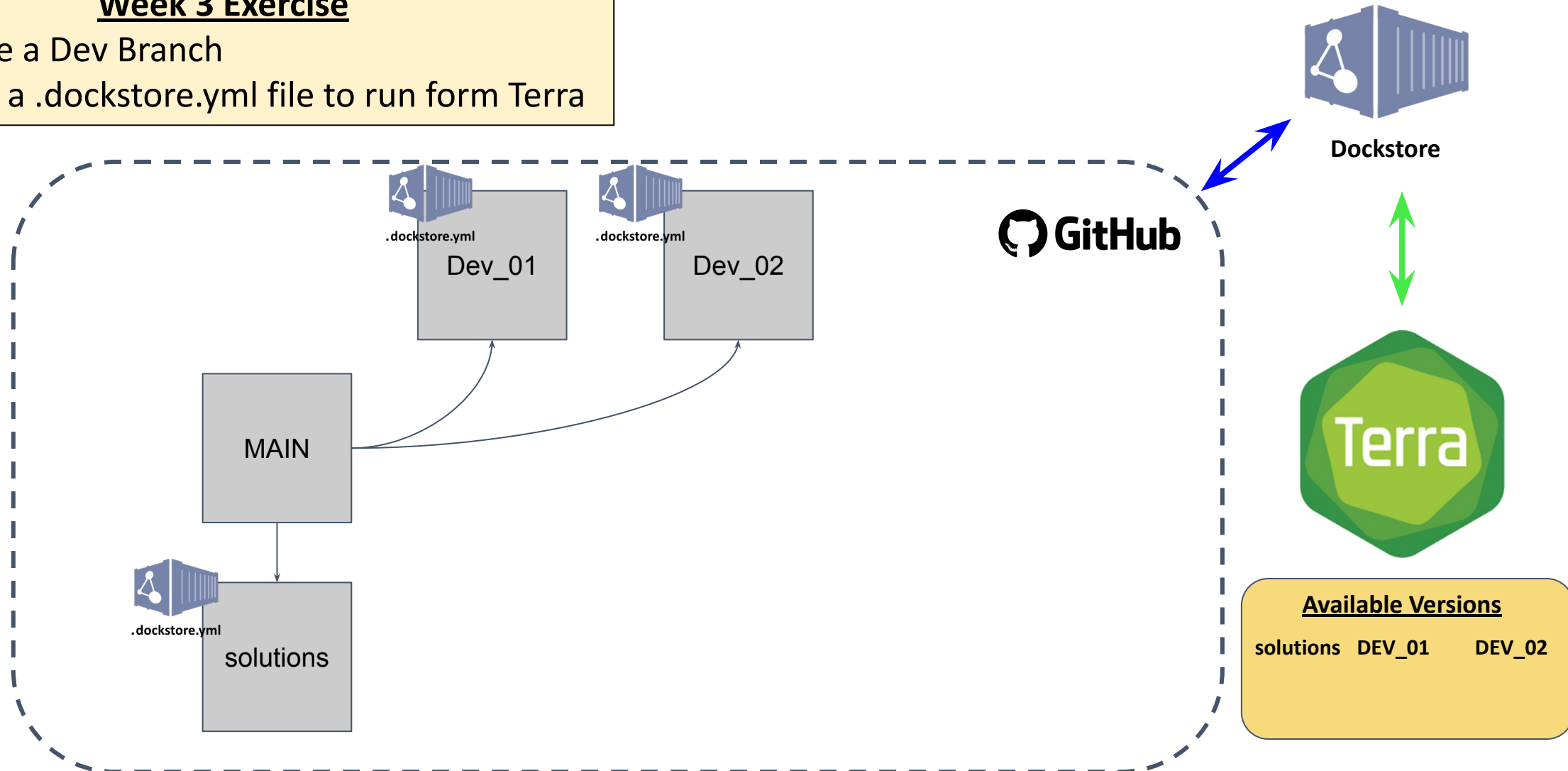- Run this workflow from your Terra workspace

# Terra Workflow Versions

# From Terra Workspace

# Terra Workflow Versions



**Week 3 Exercise**
1. Create a Dev Branch
2. Write a .dockstore.yml file to run form Terra

Dockstore

**.dockstore.yml** Dev_01

**.dockstore.yml** Dev_02

GitHub

MAIN

**.dockstore.yml** solutions

Terra

**Available Versions**

solutions   DEV_01        DEV_02

# From Terra Workspace

# Workflow Management Training

**10m to begin Exercise Part 2**

We will regroup again at 10:58AM (PT)