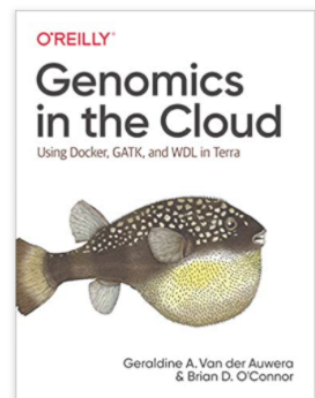# Genomics in the Cloud

## The Semi-Official Companion Booklet

**Author:** Geraldine A. Van der Auwera

**Date:** December 31, 2020

**Version:** 1.0

*Original book: https://oreil.ly/genomics-cloud*

# About this booklet

—⸮⸜⸝⸮—

This booklet contains the figures used in Genomics in the Cloud (in full color) and their captions. Its primary purpose is to provide a way for readers of the print version of the book (which is in grayscale) to access the full color versions of the figures, either by browsing the PDF or printing it out. The booklet also includes a list of chapters as well as a table of contents for each chapter, which might be helpful as a quick reference. Note that this first version is a little rough around the edges; there is a lot of opportunity for improvement, but it's going to take some wrestling with LaTeX... All feedback and offers of help are welcome!

## Figure re-use policy

The figures in this booklet are subject to the same rules/copyright as the original book. You are welcome to use any of these figures in teaching materials, presentations and so on, without express permission (except for the two figures credited to others, 6.15 and 3.3; please check their respective authors' requirements). Attribution is always appreciated, e.g. "*Genomics in the Cloud* by Geraldine A. Van der Auwera and Brian D. O'Connor (O'Reilly), https://oreil.ly/genomics-cloud".

## Additional resources

- Book: https://oreil.ly/genomics-cloud
- Blog: https://broadinstitute.github.io/genomics-in-the-cloud
- Github: https://github.com/broadinstitute/genomics-in-the-cloud
- Figures: https://console.cloud.google.com/storage/browser/genomics-in-the-cloud/figures/

## Contact information

- Twitter: @VdaGeraldine
- LinkedIn: Geraldine Van der Auwera

## Cover image credit

Brocken Inaglory, Wikimedia Commons

# List of Chapters

# List of Figures

# Chapter 1 Introduction

Why you should care about the cloud, and how bioinformatics / life sciences research benefits from moving to a cloud-based ecosystem for data sharing and analysis. No, the cloud environment is not perfect; yes, it really is a game changer.

## 1.1 The Promises and Challenges of Big Data in Biology and Life Sciences

## 1.2 Infrastructure Challenges

## 1.3 Toward a Cloud-Based Ecosystem for Data Sharing and Analysis

### 1.3.1 Cloud-Hosted Data and Compute

### 1.3.2 Platforms for Research in the Life Sciences

### 1.3.3 Standardization and Reuse of Infrastructure

## 1.4 Being FAIR

**Figure 1.1:** Recorded growth of sequencing datasets up to 2015 and projected growth for the next decade (top); growth in data production at the Broad Institute (bottom).



**Figure 1.2:** GATK provides a series of Best Practices to process sequence data for a variety of experimental designs.

**Figure 1.3:** Inverting the model for data sharing.



**Figure 1.4:** Data Biosphere principles in action: federated data analysis across multiple datasets in Terra using a workflow imported from Dockstore and executed in GCP.

# Chapter 2  Genomics in a Nutshell: A Primer for Newcomers to the Field

A primer for newcomers to the field of genomics, covering foundational terms and concepts such as genes, DNA and genomic variation, plus the technical basics of sequencing and handling genomic data.

## 2.1  Introduction to Genomics

**2.1.1**  The Gene as a Discrete Unit of Inheritance (Sort Of)

**2.1.2**  The Central Dogma of Biology: DNA to RNA to Protein

**2.1.3**  The Origins and Consequences of DNA Mutations

**2.1.4**  Genomics as an Inventory of Variation in and Among Genomes

**2.1.5**  The Challenge of Genomic Scale, by the Numbers

## 2.2  Genomic Variation

**2.2.1**  The Reference Genome as Common Framework

**2.2.2**  Physical Classification of Variants

**2.2.3**  Germline Variants Versus Somatic Alterations

## 2.3  High-Throughput Sequencing Data Generation

**2.3.1**  From Biological Sample to Huge Pile of Read Data

**2.3.2**  Types of DNA Libraries: Choosing the Right Experimental Design

## 2.4  Data Processing and Analysis

**2.4.1**  Mapping Reads to the Reference Genome

**2.4.2**  Variant Calling

**2.4.3**  Data Quality and Sources of Error

**2.4.4**  Functional Equivalence Pipeline Specification

**Figure 2.1:** The chromosome (shown here in the form of two sister chromatids, each composed of one incredibly long molecule of double-stranded DNA) on which we delineate genes composed of exons and introns.



**Figure 2.2:** The central dogma of biology: DNA leads to RNA; RNA leads to amino acids; amino acids lead to protein.

**Second letter**

|  | U | C | A | G |  |
|---|---|---|---|---|---|
| **U** | UUU UUC Phenylalanine<br>UUA UUG Leucine | UCU UCC UCA UCG Serine | UAU UAC Tyrosine<br>UAA Stop codon<br>UAG Stop codon | UGU UGC Cysteine<br>UGA Stop codon<br>UGG Tryptophan | U C A G |
| **C** | CUU CUC CUA CUG Leucine | CCU CCC CCA CCG Proline | CAU CAC Histidine<br>CAA CAG Glutamine | CGU CGC CGA CGG Arginine | U C A G |
| **A** | AUU AUC AUA Isoleucine<br>AUG Methionine; start codon | ACU ACC ACA ACG Threonine | AAU AAC Asparagine<br>AAA AAG Lysine | CGU CGC Serine<br>CGA CGG Arginine | U C A G |
| **G** | GUU GUC GUA GUG Valine | GCU GCC GCA GCG Alanine | GAU GAC Aspartic acid<br>GAA GAG Glutamic acid | GGU GGC GGA GGG Glycine | U C A G |

First letter (left) — Third letter (right)

**Figure 2.3:** The genetic code connects three-letter codons in a messenger RNA sequence to specific amino acids.

Normal gene        Mutated gene

or

Normal protein        Abnormal protein    No protein

**Figure 2.4:** A mutation in the DNA sequence can cause the gene's protein product to function abnormally or disable its production entirely.

REF    chr1                                                                 chr5

PAIRS OF READS

No coverage

Less coverage

VARIANT

Point mutation (SNP/SNV)    Indel    Homozygous deletion    Hemizygous deletion    Amplification    Translocation breakpoint

Short variants (<50 basepairs)    Copy number variants/ alterations    Structural variant

**Figure 2.5:** The major types of variant classified by physical changes to the DNA.

**Figure 2.6:** A single-nucleotide variant.



**Figure 2.7:** Indels can be insertions (left) or deletions (right).



**Figure 2.8:** Example of copy-number variant caused by a duplication.



**Figure 2.9:** Examples of structural variants.

**Figure 2.10:** Germline variants are present in all cells of the body (left) while somatic alterations are present only in a subset of cells (right).



**Figure 2.11:** Library preparation process for bulk DNA (top); alternative pathway for bulk RNA (bottom).

**Figure 2.12:** Overview of Illumina short read sequencing.



**Figure 2.13:** FASTQ and Phred scale.



**Figure 2.14:** Key elements of the SAM format: file header and read record structure.

**Figure 2.15:** The CIGAR string describes the structure of the read alignment.



**Figure 2.16:** Experimental design comparison between whole genome (top) and exome (bottom).

**Figure 2.17:** Different exome preparation kits can lead to important differences in coverage location and quantity.



**Figure 2.18:** Visual appearance of whole genome sequence (WGS, top) and exome sequence (bottom) in a genome browser.



**Figure 2.19:** Sequence divergence introduces mapping challenges and ambiguity.

**Figure 2.20:** Paired-end sequencing helps resolve mapping ambiguity.



**Figure 2.21:** Basic structure of a VCF file.



**Figure 2.22:** Pileup of reads in IGV showing several probable short variants.

**Figure 2.23:** Relative amounts of coverage provide evidence for copy-number modeling.



**Figure 2.24:** Cheat sheet of variant metrics.

**Figure 2.25:** Common sources of error in variant discovery.



**Figure 2.26:** Some biochemical properties of the DNA itself cause biases in certain regions.

# Chapter 3 Computing Technology Basics for Life Scientists

CPU, GPU, TPU, FPGA, OMG GTFO – no really, just some basic hardware terminology, plus an introduction to key concepts like parallelism, pipelining, containers and virtual machines in fairly plain language.

## 3.1 Basic Infrastructure Components and Performance Bottlenecks

**3.1.1** Types of Processor Hardware: CPU, GPU, TPU, FPGA, OMG

**3.1.2** Levels of Compute Organization: Core, Node, Cluster, and Cloud

**3.1.3** Addressing Performance Bottlenecks

## 3.2 Parallel Computing

**3.2.1** Parallelizing a Simple Analysis

**3.2.2** From Cores to Clusters and Clouds: Many Levels of Parallelism

**3.2.3** Trade-Offs of Parallelism: Speed, Efficiency, and Cost

## 3.3 Pipelining for Parallelization and Automation

**3.3.1** Workflow Languages

**3.3.2** Popular Pipelining Languages for Genomics

**3.3.3** Workflow Management Systems

## 3.4 Virtualization and the Cloud

**3.4.1** VMs and Containers

**3.4.2** Introducing the Cloud

**3.4.3** Categories of Research Use Cases for Cloud Services

**Figure 3.1:** Levels of compute organization.



**Figure 3.2:** Scatter-gather allows parallel execution of tasks on different CPU cores (on a single machine or multiple machines, depending on how it's implemented).



**Figure 3.3:** XKCD comic on the proliferation of standards (source: https://xkcd.com/927).

A. Software stack on physical machine, e.g., your laptop



B. Virtual machines



C. Containers

**Figure 3.4:** A) The software stack installed on a physical machine; B) a system hosting multiple VMs; C) a system hosting multiple containers.



**Figure 3.5:** A system with three VMs: the one on the left is running two containers, serving App 1 and App 2; the middle is running a single container, serving App 3; the right is serving App 4 directly (no container).



**Figure 3.6:** The relationship between registry, image, and container.

**Figure 3.7:** The process for creating a Docker image.

# Chapter 4  First Steps in the Cloud

Finally we get to do some hands-on work (on Google Cloud). Set up an account, get free credits, practice managing data in storage buckets and interacting with a Docker container, get a nice custom VM set up to do some genomics.

## 4.1  Setting Up Your Google Cloud Account and First Project

**4.1.1**  Creating a Project

**4.1.2**  Checking Your Billing Account and Activating Free Credits

## 4.2  Running Basic Commands in Google Cloud Shell

**4.2.1**  Logging in to the Cloud Shell VM

**4.2.2**  Using gsutil to Access and Manage Files

**4.2.3**  Pulling a Docker Image and Spinning Up the Container

**4.2.4**  Mounting a Volume to Access the Filesystem from Within the Container

## 4.3  Setting Up Your Own Custom VM

**4.3.1**  Creating and Configuring Your VM Instance

**4.3.2**  Logging into Your VM by Using SSH

**4.3.3**  Checking Your Authentication

**4.3.4**  Copying the Book Materials to Your VM

**4.3.5**  Installing Docker on Your VM

**4.3.6**  Setting Up the GATK Container Image

**4.3.7**  Stopping Your VM. . . to Stop It from Costing You Money

## 4.4  Configuring IGV to Read Data from GCS Buckets

**Figure 4.1:** Creating a new project.

**Figure 4.2:** The panel in the Billing console summarizing free trial credits availability.

**Figure 4.3:** Budget and alert threshold administration.

**Figure 4.4:** Location of the Project ID in the GCP console.



**Figure 4.5:** GCP console storage browser.



**Figure 4.6:** Naming your bucket.

**Figure 4.7:** Viewing the contents of your bucket.



**Figure 4.8:** Mounting a directory from your Google Cloud Shell VM into a Docker container: Ubuntu container used in this chapter (left); GATK container introduced in First Steps with GATK (right).

*TIP: Click the pin symbol to "pin" this service in the shortcuts menu*

**Figure 4.9:** Compute Engine menu showing the VM instances menu item.



**Figure 4.10:** Create a VM instance.

**Figure 4.11:** The VM instance configuration panel.



**Figure 4.12:** Name your VM instance.



**Figure 4.13:** Selecting a machine type.

**Figure 4.14:** Choosing a boot disk size and image.



**Figure 4.15:** Selecting a base image.



**Figure 4.16:** Setting the boot disk size.



**Figure 4.17:** The updated boot disk selection.

**Figure 4.18:** Viewing the VM status.



**Figure 4.19:** Options for SSHing into your VM.



**Figure 4.20:** VM instance terminal.



**Figure 4.21:** Stopping, starting, or deleting your VM instance.

**Figure 4.22:** Selecting the Preferences menu item.



**Figure 4.23:** The IGV Preferences pane.



**Figure 4.24:** Selecting the Google Login menu item.

**Figure 4.25:** The Load from URL menu item.



**Figure 4.26:** The Load from URL dialog box.



**Figure 4.27:** IGV view of a BAM file located in a GCS bucket.



**Figure 4.28:** Changing the behavior of the detail viewer from "on Hover" to "on Click."

# Chapter 5  First Steps with GATK

$\mathrel{-\!\!-\!\!\circ\!\ll\!\!\approx\!\!\gg\!\circ\!\!-\!\!-}$

Let's meet the workhorse of genomics! We start with a general overview, requirements, command line syntax, the usual – then dive into calling variants with HaplotypeCaller, plus some visual troubleshooting and variant filtering concepts.

## 5.1  Getting Started with GATK

**5.1.1**  Operating Requirements

**5.1.2**  Command-Line Syntax

**5.1.3**  Multithreading with Spark

**5.1.4**  Running GATK in Practice

## 5.2  Getting Started with Variant Discovery

**5.2.1**  Calling Germline SNPs and Indels with HaplotypeCaller

**5.2.2**  Filtering Based on Variant Context Annotations

## 5.3  Introducing the GATK Best Practices

**5.3.1**  Best Practices Workflows Covered in This Book

**Figure 5.1:** The four stages of HaplotypeCaller's operation.



**Figure 5.2:** The original BAM file and output VCF file loaded in IGV.

**Figure 5.3:** IGV alignment settings.



**Figure 5.4:** Turning on the display of soft clips shows a lot of information that was hidden.



**Figure 5.5:** Realigned reads in the bamout file (bottom track).

**Figure 5.6:** Bamout shows artificial haplotypes constructed by HaplotypeCaller.



**Figure 5.7:** Bamout shows support per haplotype.



**Figure 5.8:** Density plot of QUAL (left); scatter plot of QUAL versus DP (right).

**Figure 5.9:** Density plot of QUAL: all calls together (left); stratified by callsets annotation (right).



**Figure 5.10:** Density plot of QD: all calls together (left); stratified by callsets annotation (right).

**Figure 5.11:** A scatter plot with marginal densities of QD versus DP.



| | Germline | Somatic |
|---|---|---|
| **SNPs & Indels** | 📖 HaplotypeCaller / Joint Calling | 📖 MuTect2 / Tumor-Normal |
| **Copy number** | GATK gCNV | 📖 GATK CNV + aCNV |
| **Structural Variation** | GATK SVDiscovery (beta) | *on the roadmap* |

**Figure 5.12:** Table of standard variant discovery use cases covered by GATK Best Practices.

# Chapter 6  GATK Best Practices for Germline Short Variant Discovery

Step by step examination of what may be the most commonly run genomics pipeline in the world, with highlights on joint calling for populations and deep learning for single-sample analysis.

## 6.1  Data Preprocessing

**6.1.1**  Mapping Reads to the Genome Reference

**6.1.2**  Marking Duplicates

**6.1.3**  Recalibrating Base Quality Scores

## 6.2  Joint Discovery Analysis

**6.2.1**  Overview of the Joint Calling Workflow

**6.2.2**  Calling Variants per Sample to Generate GVCFs

**6.2.3**  Consolidating GVCFs

**6.2.4**  Applying Joint Genotyping to Multiple Samples

**6.2.5**  Filtering the Joint Callset with Variant Quality Score Recalibration

**6.2.6**  Refining Genotype Assignments and Adjusting Genotype Confidence

**6.2.7**  Next Steps and Further Reading

## 6.3  Single-Sample Calling with CNN Filtering

**6.3.1**  Overview of the CNN Single-Sample Workflow

**6.3.2**  Applying 1D CNN to Filter a Single-Sample WGS Callset

**6.3.3**  Applying 2D CNN to Include Read Data in the Modeling

**Figure 6.1:** The main steps in the preprocessing workflow.

**Figure 6.2:** Reads marked as duplicates because they originated from the same DNA fragment in the library.



**Figure 6.3:** The effect of duplicate marking visualized in Integrated Genome Viewer.

**Figure 6.4:** Visualizing the effect of BQSR.



**Figure 6.5:** Sites that would be omitted from the VCF in a single-sample callset.

**Figure 6.6:** Seeing concordant evidence in multiple samples boosts our confidence that there is real variation.



**Figure 6.7:** Traditional multisample analysis scales poorly and causes the N + 1 problem.

**Figure 6.8:** The GVCF workflow improves the scaling of joint calling and solves the N + 1 problem.



**Figure 6.9:** Progression from per-sample GVCFs to final cohort VCF.



**Figure 6.10:** GVCFs viewed in IGV show tiled nonvariant blocks.

**Figure 6.11:** Variant call with genotype assignment for the three samples.



**Figure 6.12:** Gaussian clusters learned from a training set are applied to novel variant calls.



**Figure 6.13:** How the VQSLOD score is calculated for an individual annotation.



**Figure 6.14:** Genotype assignments corrected on the basis of pedigree and population priors.

**Figure 6.15:** Labradoodle or fried chicken? (Source: Karen Zack, @teenybiscuit).



**Figure 6.16:** Different calls made by 1D and 2D CNN models.

# Chapter 7  GATK Best Practices for Somatic Variant Discovery

Switching gears to cancer genomics with a rundown of how somatic calling is different; step by step through the pipelines for somatic short variants (Mutect2) and copy number alterations.

## 7.1  Challenges in Cancer Genomics

## 7.2  Somatic Short Variants (SNVs and Indels)

**7.2.1** Overview of the Tumor-Normal Pair Analysis Workflow

**7.2.2** Creating a Mutect2 PoN

**7.2.3** Running Mutect2 on the Tumor-Normal Pair

**7.2.4** Estimating Cross-Sample Contamination

**7.2.5** Filtering Mutect2 Calls

**7.2.6** Annotating Predicted Functional Effects with Funcotator

## 7.3  Somatic Copy-Number Alterations

**7.3.1** Overview of the Tumor-Only Analysis Workflow

**7.3.2** Collecting Coverage Counts

**7.3.3** Creating a Somatic CNA PoN

**7.3.4** Applying Denoising

**7.3.5** Performing Segmentation and Call CNAs

**7.3.6** Additional Analysis Options

**Figure 7.1:** Tumor progression leads to heterogeneity (left); sampling is difficult (right).



**Figure 7.2:** The fundamental concept of Tumor-Normal comparison.

**Figure 7.3:** Best Practices for somatic short variant discovery.



**Figure 7.4:** Zooming in on TP53 in IGV.

**Figure 7.5:** Difference between copy number and copy ratio.



**Figure 7.6:** Spectral karyotyping paints each chromosome pair with a color, showing various chromosomal segments that are amplified or missing (colors in left and right panels are not expected to match).

**Figure 7.7:** Best Practices workflow for somatic copy-number alteration discovery.

**Figure 7.8:** Read counts in each genomic target or bin form the basis for estimating segmented copy ratio, and each dot is the value for a single target or bin.



**Figure 7.9:** Copy-number alteration analysis plots showing the standardized copy ratios after the first step of denoising (top) and the fully denoised copy ratios after the second round (bottom).



**Figure 7.10:** Plot of segments modeled based on denoised copy ratios.

**Figure 7.11:** Full progression from raw data to results.

# Chapter 8 Automating Analysis Execution with Workflows

Halfway point; we pivot to the challenges of automating and scaling up these analyses, introducing the Cromwell workflow system and the portable Workflow Description Language (WDL).

## 8.1 Introducing WDL and Cromwell

## 8.2 Installing and Setting Up Cromwell

## 8.3 Your First WDL: Hello World

**8.3.1** Learning Basic WDL Syntax Through a Minimalist Example

**8.3.2** Running a Simple WDL with Cromwell on Your Google VM

**8.3.3** Interpreting the Important Parts of Cromwell's Logging Output

**8.3.4** Adding a Variable and Providing Inputs via JSON

**8.3.5** Adding Another Task to Make It a Proper Workflow

## 8.4 Your First GATK Workflow: Hello HaplotypeCaller

**8.4.1** Exploring the WDL

**8.4.2** Generating the Inputs JSON

**8.4.3** Running the Workflow

**8.4.4** Breaking the Workflow to Test Syntax Validation and Error Messaging

## 8.5 Introducing Scatter-Gather Parallelism

**8.5.1** Exploring the WDL

**8.5.2** Generating a Graph Diagram for Visualization

**Figure 8.1:** A hypothetical workflow that runs HaplotypeCaller.



**Figure 8.2:** A workflow that parallelizes the execution of HaplotypeCaller.

**Figure 8.3:** Visualizing the workflow graph in an online Graphviz application.

# Chapter 9  Deciphering Real Genomics Workflows

—◦◦◦◦◦—

We pretend to stumble across 2 mystery workflows, go through a systematic process of investigating their content to understand what they do and how they do it, learning useful WDL features along the way.

## 9.1  Mystery Workflow #1: Flexibility Through Conditionals

**9.1.1**  Mapping Out the Workflow

**9.1.2**  Reverse Engineering the Conditional Switch

## 9.2  Mystery Workflow #2: Modularity and Code Reuse

**9.2.1**  Mapping Out the Workflow

**9.2.2**  Unpacking the Nesting Dolls

**Figure 9.1:** Graph description in JSON (left) and visual rendering (right).



**Figure 9.2:** Visual rendering of the workflow graph.



**Figure 9.3:** Graph diagram of the VariantCalling.wdl workflow.

# Chapter 10  Running Single Workflows at Scale with Pipelines API

So far we've been running everything on our little custom VM. Now it's time to unleash the full power of the cloud by dispatching workflow tasks to multiple machines – with surprisingly little effort.

## 10.1  Introducing the GCP Genomics Pipelines API Service

**10.1.1**  Enabling Genomics API and Related APIs in Your Google Cloud Project

## 10.2  Directly Dispatching Cromwell Jobs to PAPI

**10.2.1**  Configuring Cromwell to Communicate with PAPI

**10.2.2**  Running Scattered HaplotypeCaller via PAPI

**10.2.3**  Monitoring Workflow Execution on Google Compute Engine

## 10.3  Understanding and Optimizing Workflow Efficiency

**10.3.1**  Granularity of Operations

**10.3.2**  Balance of Time Versus Money

**10.3.3**  Suggested Cost-Saving Optimizations

**10.3.4**  Platform-Specific Optimization Versus Portability

## 10.4  Wrapping Cromwell and PAPI Execution with WDL Runner

**10.4.1**  Setting Up WDL Runner

**10.4.2**  Running the Scattered HaplotypeCaller Workflow with WDL Runner

**10.4.3**  Monitoring WDL Runner Execution

**Figure 10.1:** Overview of Cromwell + PAPI operation.



**Figure 10.2:** Logos and descriptions for the three required APIs: Genomics API, Cloud Storage JSON API, and Compute Engine API.



**Figure 10.3:** Side-by-side comparison of local versus PAPI execution.

**Figure 10.4:** List of active VM instances.



**Figure 10.5:** Overview of Compute Engine activity.

**Figure 10.6:** Overview of WDL Runner operation.



**Figure 10.7:** List of active VM instances (WDL Runner submission).



**Figure 10.8:** Output from the WDL Runner submission.

# Chapter 11  Running Many Workflows Conveniently in Terra

—◦◦◦◦◦—

Now we're scaling up to arbitrary numbers of samples, using the managed Cromwell server in Terra, an open platform for secure data access and analysis.

## 11.1  Getting Started with Terra

**11.1.1**  Creating an Account

**11.1.2**  Creating a Billing Project

**11.1.3**  Cloning the Preconfigured Workspace

## 11.2  Running Workflows with the Cromwell Server in Terra

**11.2.1**  Running a Workflow on a Single Sample

**11.2.2**  Running a Workflow on Multiple Samples in a Data Table

**11.2.3**  Monitoring Workflow Execution

**11.2.4**  Locating Workflow Outputs in the Data Table

**11.2.5**  Running the Same Workflow Again to Demonstrate Call Caching

## 11.3  Running a Real GATK Best Practices Pipeline at Full Scale

**11.3.1**  Finding and Cloning the GATK Best Practices Workspace for Germline Short Variant Discovery

**11.3.2**  Examining the Preloaded Data

**11.3.3**  Selecting Data and Configuring the Full-Scale Workflow

**11.3.4**  Launching the Full-Scale Workflow and Monitoring Execution

**11.3.5**  Options for Downloading Output Data—or Not

**Figure 11.1:** Overview of the Terra platform.

**Figure 11.2:** Expanded side menu showing sign-in button.

**Figure 11.3:** The New User Registration form.



**Figure 11.4:** The GCP console Billing account permissions panel.

**Figure 11.5:** Adding the Terra billing user account as a user on a GCP billing account.

## Add members to "My Billing Account"

### Add members and roles for "My Billing Account" resource

Enter one or more members below. Then select a role for these members to grant them access to your resources. Multiple roles allowed. Learn more

**New members**

terra-billing@terra.bio ✕

**Role**

≡  Type to filter

| Billing | Billing Account Administrator |
| Cloud Composer | Billing Account User |
| Dataflow | Billing Account Viewer |
| Dataproc | |
| Error Reporting | |
| Firebase | |
| IAM | |
| Logging | |

**MANAGE ROLES**

**Figure 11.6:** Using an existing billing account to create a billing project in Terra.

Figure 11.7: Cloning the preconfigured workspace. A) List of available actions; B) cloning form.



Figure 11.8: List of available workflow configurations.



Figure 11.9: Viewing the workflow information summary.

⊙ scatter-hc.filepaths

Snapshot:  1  ⌄

Source: genomics-in-the-cloud/scatter-hc/1
Synopsis: Run GATK4 HaplotypeCaller parallelized by interval

⌄ This workflow runs the HaplotypeCaller tool from GATK4 in GVCF mode on a single sample in BAM format. The execution of the HaplotypeCaller tool is parallelized using an intervals list file. The per-interval output GVCF files are then merged to produce a single GVCF file for the sample, which can then be used by the joint-discovery workflow according to the GATK Best Practices for germline short variant discovery.

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

**Figure 11.10:** Viewing the workflow script.

```
SCRIPT    • •    INPUTS    • •    OUTPUTS    • •    RUN ANALYSIS

1    ## This workflow runs the HaplotypeCaller tool from GATK4 in GVCF mode
2    ## on a single sample in BAM format. The execution of the HaplotypeCaller
3    ## tool is parallelized using an intervals list file. The per-interval
4    ## output GVCF files are then merged to produce a single GVCF file for
5    ## the sample, which can then be used by the joint-discovery workflow
6    ## according to the GATK Best Practices for germline short variant
7    ## discovery.
8
9    version 1.0
10
11   workflow ScatterHaplotypeCallerGVCF {
12
13       input {
14           File input_bam
15           File input_bam_index
16           File intervals_list
17       }
18
19       String output_basename = basename(input_bam, ".bam")
20
21       Array[String] calling_intervals = read_lines(intervals_list)
22
23       scatter(interval in calling_intervals) {
24           call HaplotypeCallerGVCF {
```

**Figure 11.11:** Viewing the workflow inputs.

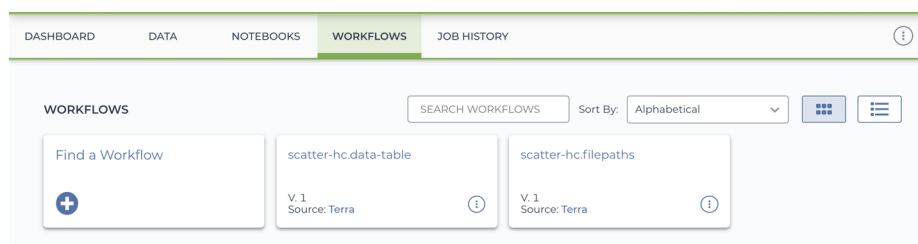| Task name | ↓ | Variable | Type | Attribute | |
|---|---|---|---|---|---|
| HaplotypeCallerGVCF | | docker_image | String | "us.gcr.io/broad-gatk/gatk:4.1.3.0" | [...] |
| HaplotypeCallerGVCF | | java_opt | String | "-Xmx8G" | [...] |
| HaplotypeCallerGVCF | | ref_dict | File | "gs://genomics-in-the-cloud/v1/data/germline/ref/ref.dict" | [...] |
| HaplotypeCallerGVCF | | ref_fasta | File | "gs://genomics-in-the-cloud/v1/data/germline/ref/ref.fasta" | [...] |

Download json | Drag or click to upload json    SEARCH INPUTS

**Figure 11.12:** The workflow launch dialog.

## Confirm launch

This analysis will be run by Cromwell 49.

This will launch **1** analysis.

CANCEL        LAUNCH

**Figure 11.13:** Overview of workflow submission in Terra.

**Figure 11.14:** The second workflow is set to run on rows in a data table.



**Figure 11.15:** The workflow input configuration references data tables.



**Figure 11.16:** Viewing the menu of data tables on the DATA tab.

**Figure 11.17:** The Workspace Data table.

| Key | Value |
| --- | --- |
| gatk_docker | us.gcr.io/broad-gatk/gatk:4.1.3.0 |
| intervals_list_full | snippet-intervals-full.list |
| intervals_list_min | snippet-intervals-min.list |
| ref_dict | ref.dict |
| ref_fasta | ref.fasta |
| ref_fasta_index | ref.fasta.fai |

**Figure 11.18:** The $book_sample table$.

| ☐ ▾ | book_sample_id ↓ | input_bam | input_bam_index |
| --- | --- | --- | --- |
| ☐ | father | father.bam | father.bai |
| ☐ | mother | mother.bam | mother.bai |
| ☐ | son | son.bam | son.bai |

**Figure 11.19:** Initiating an analysis directly on a subset of data.

**Figure 11.20:** Specifying a workflow to run on the selected data.



**Figure 11.21:** Configuration updated with data selection.



**Figure 11.22:** List of submissions in the Job History.



**Figure 11.23:** The workflow submission summary page.



**Figure 11.24:** Workflow in A) Running state and, B) Succeeded state.

A.

B.

**Figure 11.25:** A workflow in Failed state with ERRORS summary and Failure Message.



**Figure 11.26:** List of tasks and related resources.



**Figure 11.27:** Viewing the status of shards for a scattered task.



**Figure 11.28:** A timing diagram showing the breakdown of runtime per stage of execution for each task call.

**Figure 11.29:** A timing diagram showing preempted calls (green bars, at lines 2, 12, and 13 from the top).



**Figure 11.30:** The data table showing the newly generated output$_g vcf column$.



**Figure 11.31:** The workflow outputs configuration panel.



**Figure 11.32:** The file browser interface showing workflow outputs in the workspace bucket.

**Figure 11.33:** A timing diagram showing CallCacheReading stage run time.



**Figure 11.34:** Overview of Cromwell's call caching mechanism..



**Figure 11.35:** Summary information for the Whole-Genome-Analysis-Pipeline workspace.



**Figure 11.36:** A list of tables and detailed view of the sample table.

| | sample_id | ↓ | flowcell_unmapped_bams_list | output_bqsr_reports | output_cram |
|---|---|---|---|---|---|
| ☐ | NA12878 | | NA12878.ubams.list | NA12878.recal_data.csv | NA12878.cram |
| ☐ | NA12878_small | | NA12878_24RG_small.txt | NA12878_small.recal_data.csv | NA12878_small.cram |

**TABLES** ⊕

☰ participant (1)

☰ sample (2)

**REFERENCE DATA** ⊕

☰ hg38 ⊖

⬇ DOWNLOAD ALL ROWS    ⧉ COPY PAGE TO CLIPBOARD    | 0 rows selected

**Figure 11.37:** The List View of the task calls in the master workflow.



**Figure 11.38:** The timing diagram for the master workflow showing subworkflows (solid red bars) and individual tasks that are not bundled into subworkflows (multicolor bars).

**Figure 11.39:** The workflow details page for the BamToGvcf subworkflow.



**Figure 11.40:** File download windows showing A) the list of unmapped BAM files, and B) the final GVCF output.

# Chapter 12  Interactive Analysis in Jupyter Notebook

Circling back to the GATK work from earlier chapters, we examine what that would all look like done in Jupyter Notebooks instead of the terminal shell. Between embedded IGV and ggplots galore, it looks good!

## 1.1 Hello Python

Let's try a basic Hello World example in Python.

```
In [1]:  print ("Hello World")

         Hello World

In [ ]:  # Now you try adding a variable
         greeting =
```

Figure 12.1: Doc text, code cell, and execution output in a Jupyter notebook.



Figure 12.2: An overview of the Jupyter service in Terra.



Figure 12.3: Options for customizing the software installed in the notebook runtime.

**Figure 12.4:** Notebooks in shared workspaces are protected from overwriting when two people open them concurrently.



**Figure 12.5:** The Notebooks tab showing two copies of the notebook: one already executed and another without any previous results.



**Figure 12.6:** The Notebook Runtime status widget.

**RUNTIME CONFIGURATION**                                                              ✕

Create a cloud compute instance to launch Jupyter Notebooks or a Project-Specific software application.

**ENVIRONMENT** ⓘ

New Default (released on January 14): (GATK 4.1.4.1, Python 3.7.6, R 3.6.2)        ⌄

What's installed on this environment?                          Updated: Feb 25, 2020
                                                               Version: 0.0.13

**COMPUTE POWER**

Select from one of the default runtime profiles or define your own

Profile          Default (Moderate) computer power                              ⌄

CPUs        4          Memory (GB)    15       Disk size (GB)    50

       COST: $0.19 per hour

DELETE RUNTIME                                        CANCEL        REPLACE

**Figure 12.7:** The default Notebook Runtime configuration settings.

**INSTALLED PACKAGES**                                                    ←    ✕

New Default (released on January 14): (GATK 4.1.4.1, Python 3.7.6, R 3.6.2)        ⌄

Updated: Feb 25, 2020
Version: 0.0.13

**Installed packages**    Python        ⌄

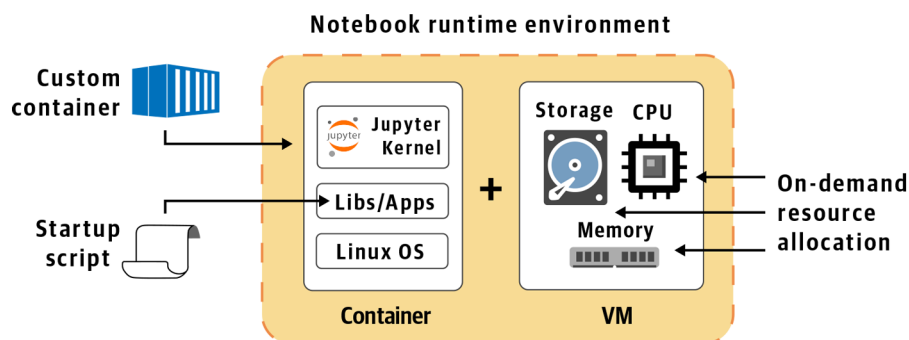| **Package** | Python ✓ | **Version** |
|---|---|---|
| lazy-object-proxy | R | 1.4.3 |
| pandocfilters | | 1.4.2 |
| googleapis-commo | Tools | 1.51.0 |
| biopython | | 1.72 |
| tf-estimator-nightly | | 1.14.0.dev2019030115 |
| ipython-genutils | | 0.2.0 |

**Figure 12.8:** Detailed view of the packages installed on the default runtime environment.

**Figure 12.9:** The Compute Power section allows you to specify a startup script if you choose the Custom profile.



**Figure 12.10:** Menu on the notebook preview page displaying the main options: Preview, Edit, and Playground Mode.



**Figure 12.11:** The standard Jupyter menu bar.



**Figure 12.12:** A newly created IGV browser.

**Figure 12.13:** The IGV browser showing the two sequence data tracks.



**Figure 12.14:** IGV.js rendering of the sequencing data ("Mother WGS" track) and output variants produced by HaplotypeCaller ("Mother variants" track).

**Figure 12.15:** Menu of display options for the Mother WGS sequence data track.

**Figure 12.16:** Display of soft clips.



**Figure 12.17:** QUAL distribution.

**Figure 12.18:** QUAL density plot.



**Figure 12.19:** QUAL density plots by callsets from GiaB.

**Figure 12.20:** Scatter plot QUAL versus DP.



**Figure 12.21:** A scatter plot along with density plots.

# Chapter 13 Assembling Your Own Workspace in Terra

——⊸∘≪∞≫∘⊸——

Crossing the bridge from canned examples to importing your own data and methods into Terra in a few different scenarios. Draws on other services in the ecosystem including Dockstore and data repositories.

## 13.1 Managing Data Inside and Outside of Workspaces

**13.1.1** The Workspace Bucket as Data Repository

**13.1.2** Accessing Private Data That You Manage Outside of Terra

**13.1.3** Accessing Data in the Terra Data Library

## 13.2 Re-Creating the Tutorial Workspace from Base Components

**13.2.1** Creating a New Workspace
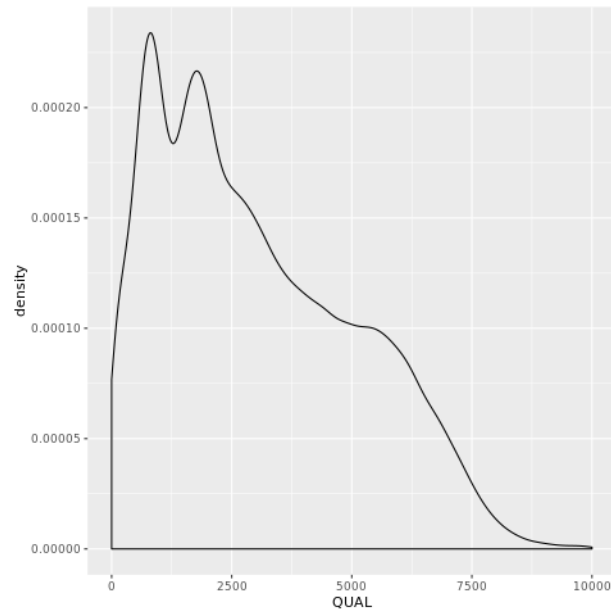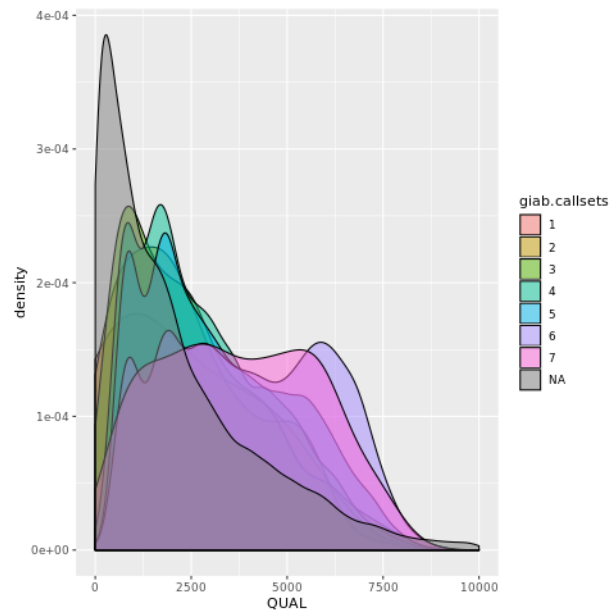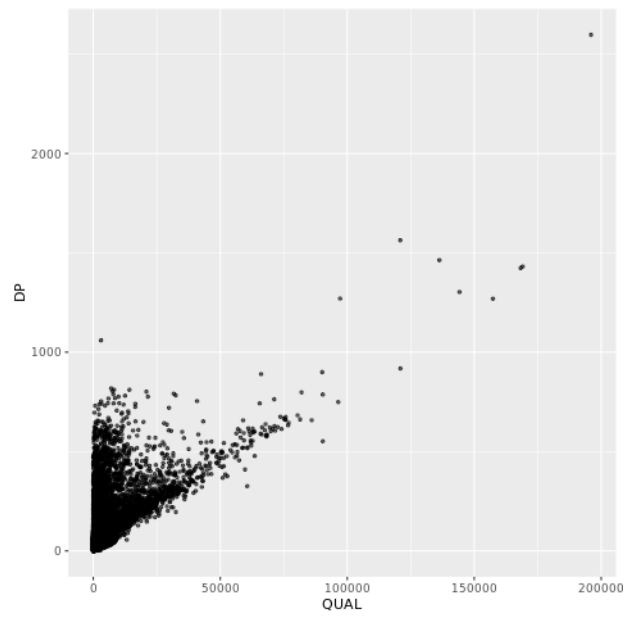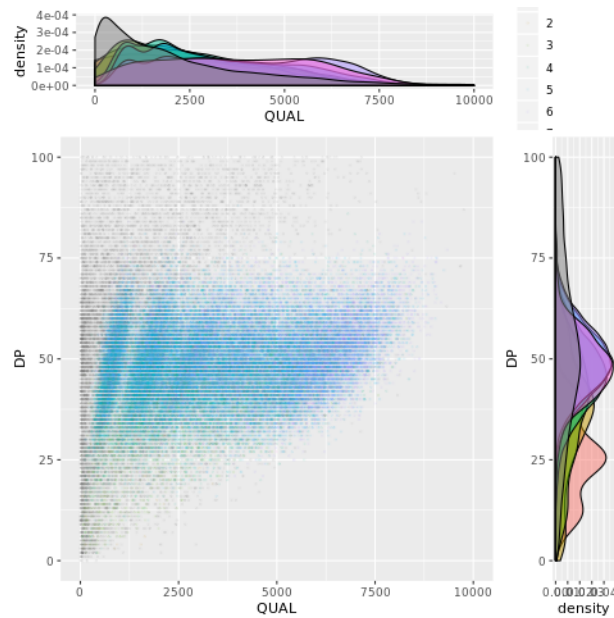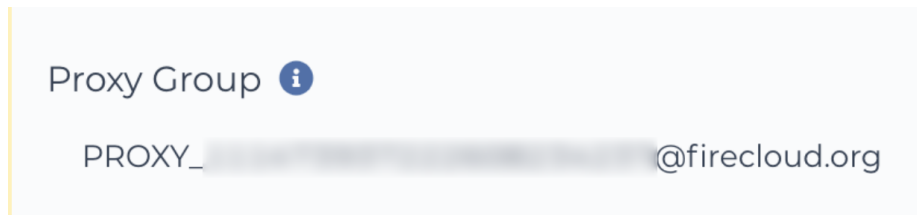
**13.2.2** Adding the Workflow to the Methods Repository and Importing It into the Workspace

**13.2.3** Creating a Configuration Quickly with a JSON File

**13.2.4** Adding the Data Table

**13.2.5** Filling in the Workspace Resource Data Table

**13.2.6** Creating a Workflow Configuration That Uses the Data Tables

**13.2.7** Adding the Notebook and Checking the Runtime Environment

**13.2.8** Documenting Your Workspace and Sharing It

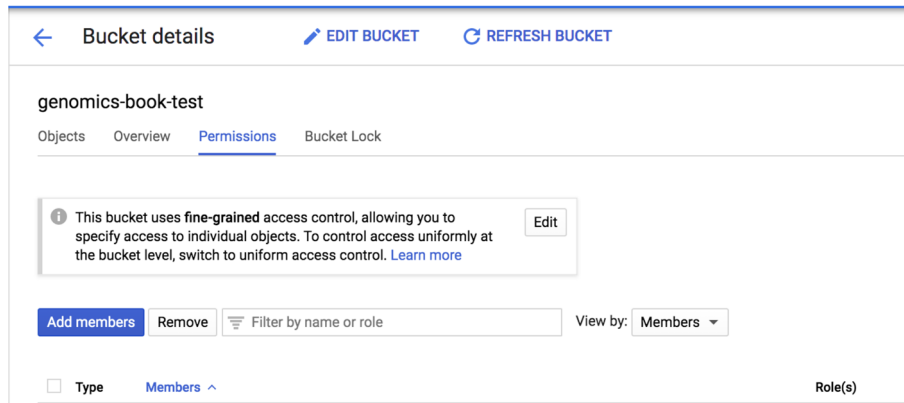## 13.3 Starting from a GATK Best Practices Workspace

**13.3.1** Cloning a GATK Best Practices Workspace

**13.3.2** Examining GATK Workspace Data Tables to Understand How the Data Is Structured

**13.3.3** Getting to Know the 1000 Genomes High Coverage Dataset

**13.3.4** Copying Data Tables from the 1000 Genomes Workspace

**13.3.5** Using TSV Load Files to Import Data from the 1000 Genomes Workspace

**13.3.6** Running a Joint-Calling Analysis on the Federated Dataset

## 13.4 Building a Workspace Around a Dataset

**13.4.1** Cloning the 1000 Genomes Data Workspace

**13.4.2** Importing a Workflow from Dockstore

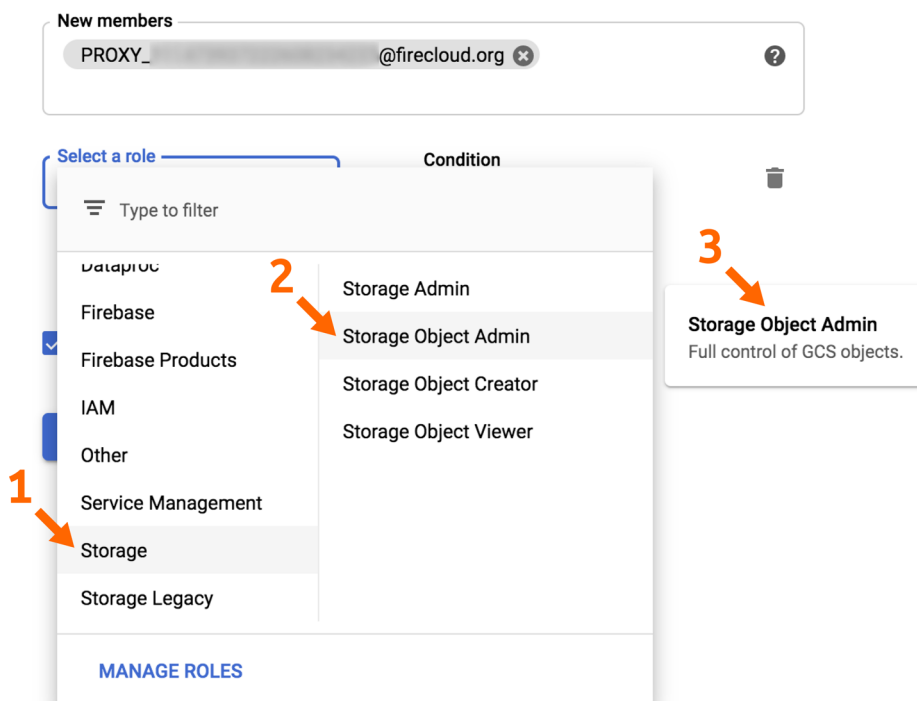**13.4.3** Configuring the Workflow to Use the Data Tables

**Figure 13.1:** The proxy group identifier displayed in the user profile.



**Figure 13.2:** The bucket permissions panel showing accounts with access to the bucket.

**Figure 13.3:** Granting access to a bucket to a new member.
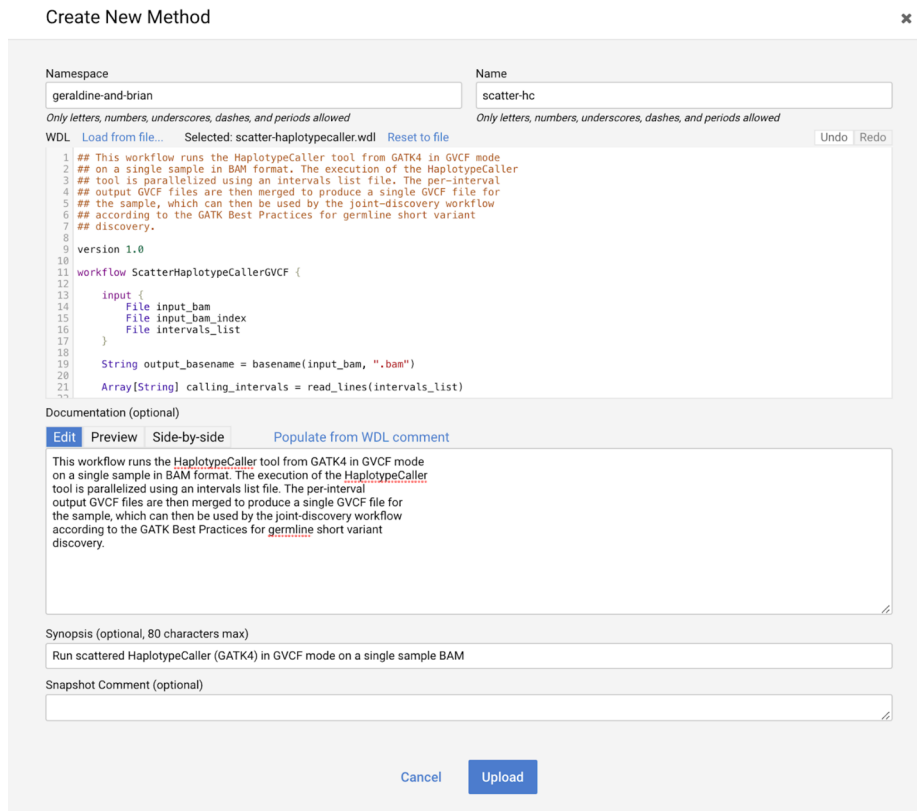
**Figure 13.4:** The Create a New Workspace dialog box.

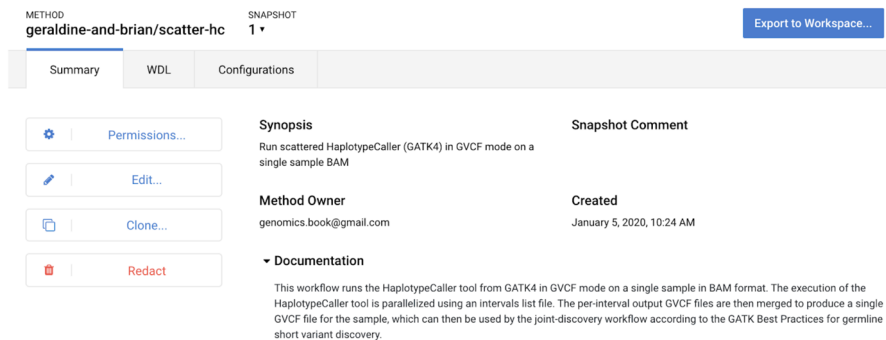**Figure 13.5:** The Create New Method page in the Broad Methods Repository.



**Figure 13.6:** Summary page for the newly created workflow.



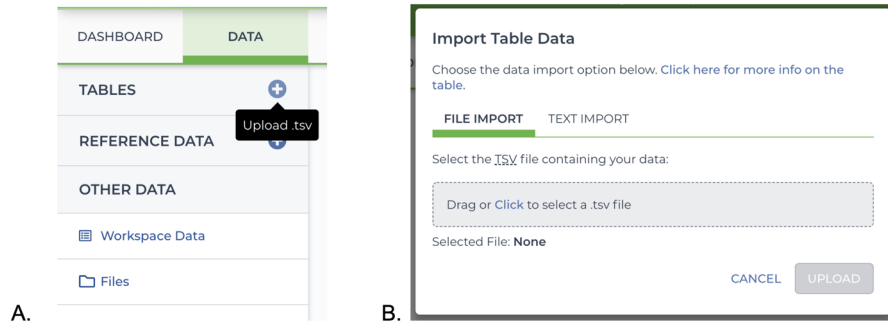**Figure 13.7:** A sample data table from the tutorial workspace, viewed in Google Sheets.

**Figure 13.8:** TSV load file import A) button, and B) dialog.



**Figure 13.9:** The data model—the structure of the example dataset.
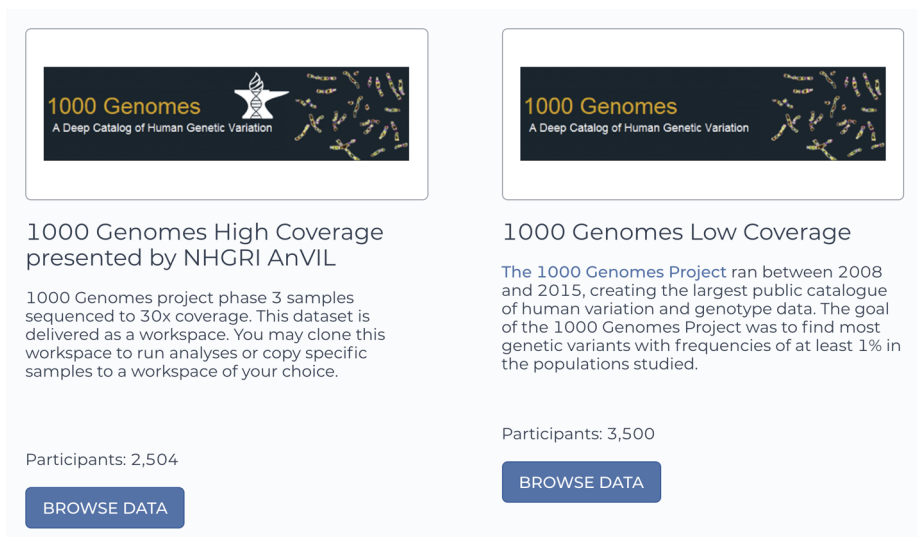


**Figure 13.10:** The Terra Data Library contains two repositories of data from the 1000 Genomes Project.



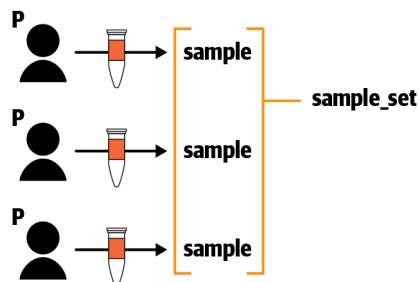**Figure 13.11:** The data model for the 1000 Genomes High Coverage dataset.
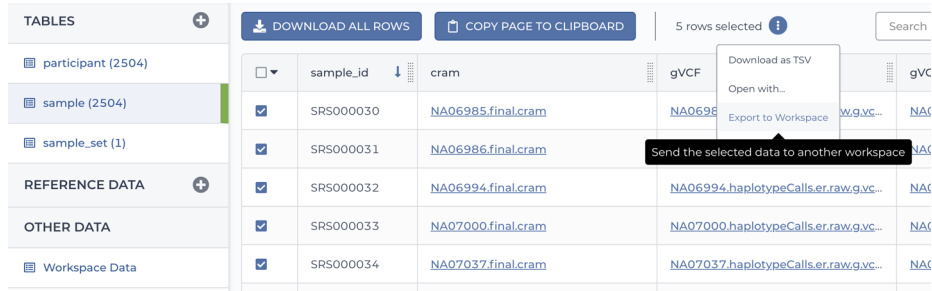
**Figure 13.12:** The Copy Data to Workspace dialog box.



**Figure 13.13:** Direct text import of TSV-formatted data table content.



...



**Figure 13.14:** Start and end rows of the membership load file $sample_set_membership.tsv$.

Figure 13.15: Updated membership load file $sample_set_membership.tsv$ $assigning 25 samples to the federated-dataset sample set$.



Figure 13.16: The $sample_set$ $table showing the three sample sets$.



Figure 13.17: Input configuration details for the $input_gvcfs$ and $input_gvcfs_indices$ variables.
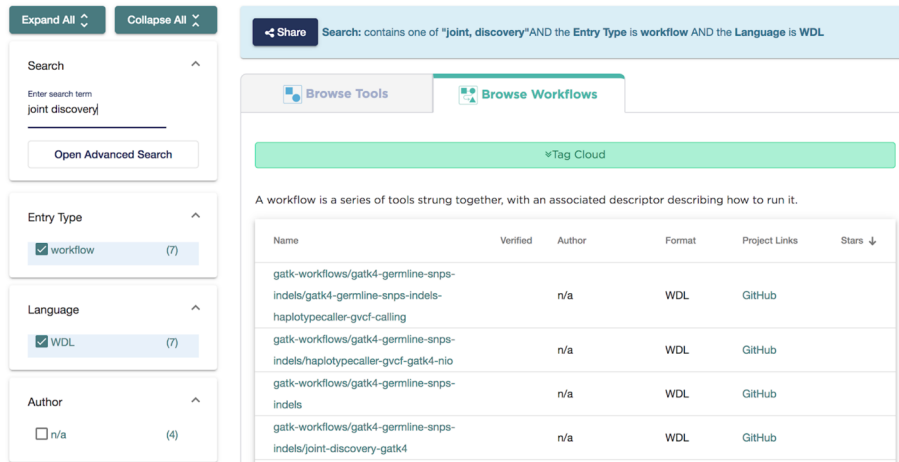
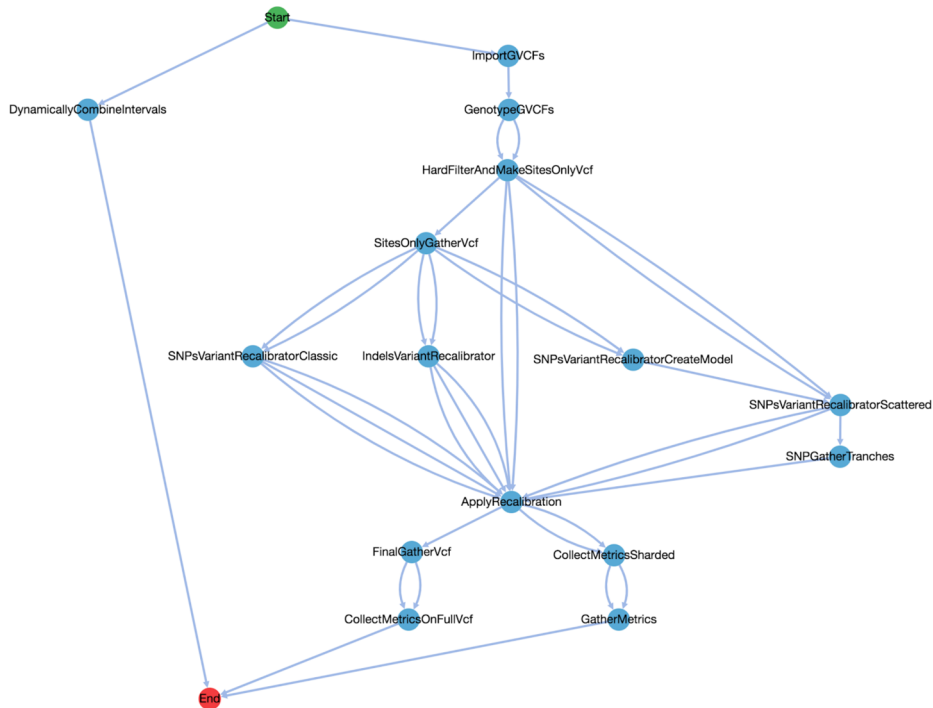**Figure 13.18:** Search results for "joint discovery" in Dockstore.



**Figure 13.19:** The Joint Discovery workflow provided in the DAG tab in Dockstore.

# Chapter 14  Making a Fully Reproducible Paper

— ⋘∘◊∘⋙ —

Capstone case study on computational reproducibility involving synthetic data creation, GATK, downstream analysis and real biological findings by Dr. Matthieu Miossec et al.

## 14.1  Overview of the Case Study

**14.1.1**  Computational Reproducibility and the FAIR Framework

**14.1.2**  Original Research Study and History of the Case Study

**14.1.3**  Assessing the Available Information and Key Challenges

**14.1.4**  Designing a Reproducible Implementation

## 14.2  Generating a Synthetic Dataset as a Stand-In for the Private Data

**14.2.1**  Overall Methodology

**14.2.2**  Retrieving the Variant Data from 1000 Genomes Participants

**14.2.3**  Creating Fake Exomes Based on Real People

**14.2.4**  Mutating the Fake Exomes

**14.2.5**  Generating the Definitive Dataset

## 14.3  Re-Creating the Data Processing and Analysis Methodology

**14.3.1**  Mapping and Variant Discovery

**14.3.2**  Variant Effect Prediction, Prioritization, and Variant Load Analysis

**14.3.3**  Analytical Performance of the New Implementation

## 14.4  The Long, Winding Road to FAIRness
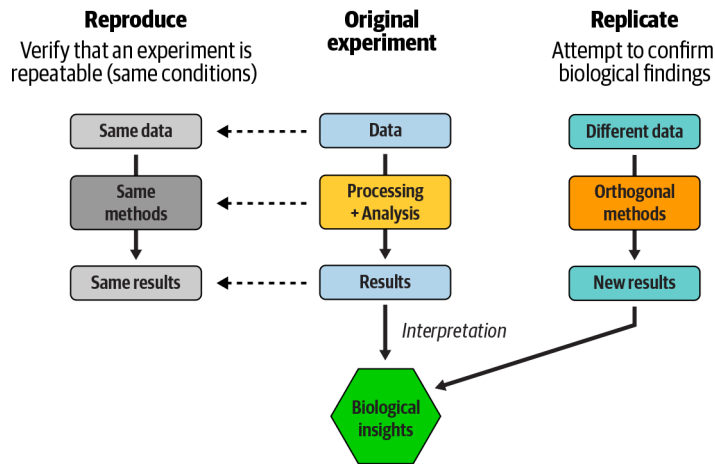
## 14.5  Final Conclusions

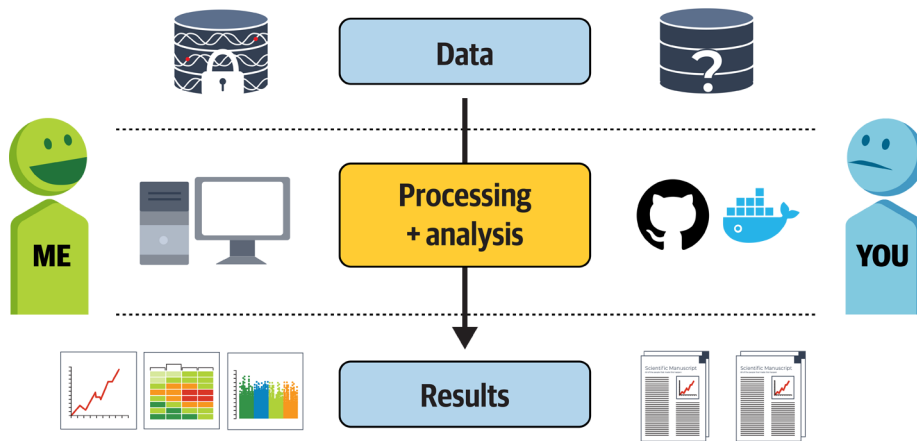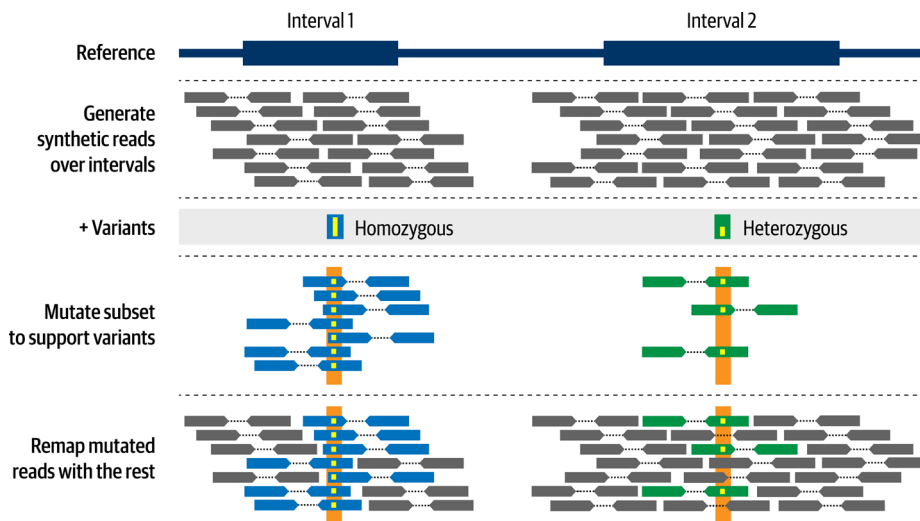**Figure 14.1:** Reproducibility of an analysis versus replicability of study findings.



**Figure 14.2:** Typical asymmetry in the availability of information between author and reader.
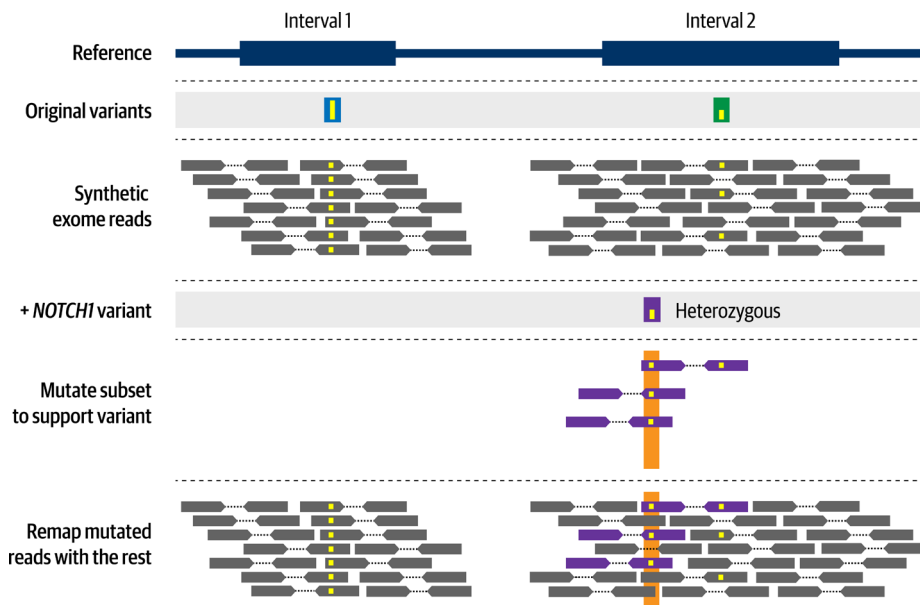


**Figure 14.3:** Summary of the information provided in the original preprint of the Tetralogy of Fallot paper.

**Figure 14.4:** Replacing a real dataset that cannot be distributed with a synthetic dataset that mimics the original data's characteristics.



**Figure 14.5:** Overview of our implementation for generating appropriate synthetic data.
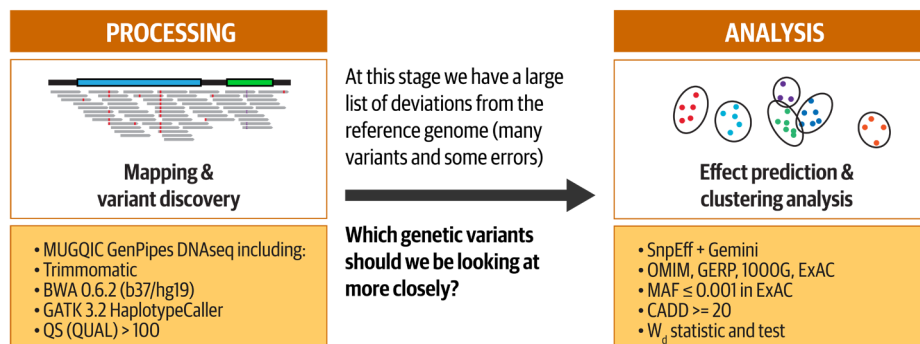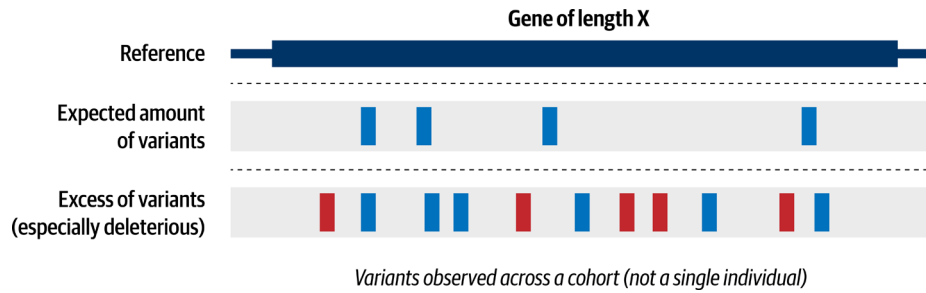
**Figure 14.6:** NEAT-genReads creates simulated read data based on a reference genome and list of variants.
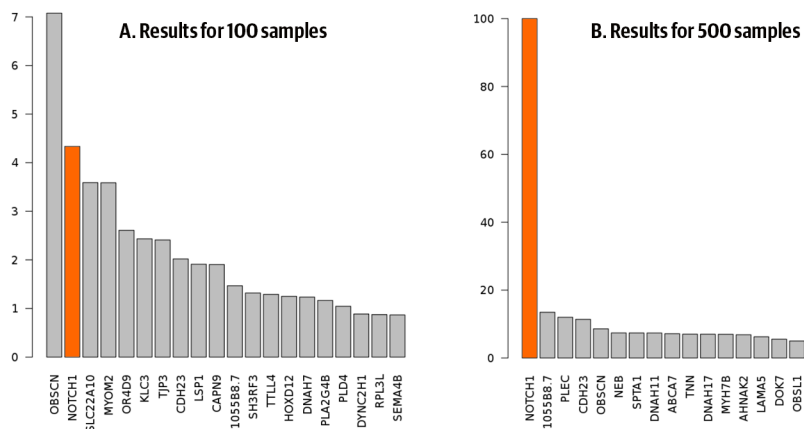


**Figure 14.7:** BAMSurgeon introduces mutations in read data.



**Figure 14.8:** Summary of the two phases of the study: Processing and Analysis.

**Figure 14.9:** Comparing variant load in a gene across multiple samples.



**Figure 14.10:** Ranking from the clustering test for A) 100-participant set, and B) 500-participant set.

# End notes

2020 has been a rough year.

Let's all work together to make 2021 more safe, equitable and enjoyable for all.

Best wishes and don't hesitate to ask for help!