# What does reproducibility mean to you?

Reproduce ≠ Replicate

OTHER DATA

OTHER ANALYSIS

"Can we confirm that this is true?"

DATA → ANALYSIS → RESULT → BIOLOGICAL FINDINGS

? → ANALYSIS → RESULT

"Can we repeat the experiment?"

The original project broad.io/ASHG2018

This hackathon project broad.io/FAIRdatahack2019

1-Collect-1000G-participant

Variant calls

2-Generate-synthetic-reads

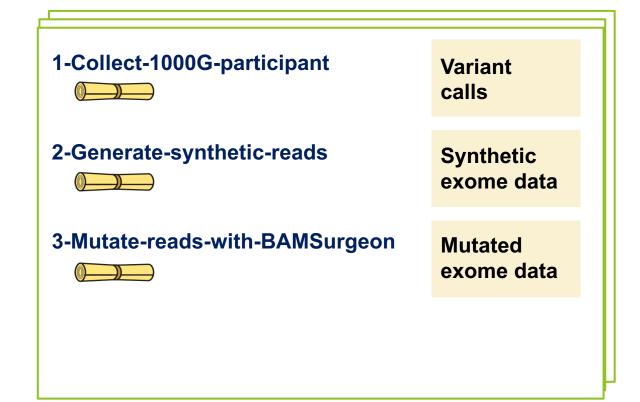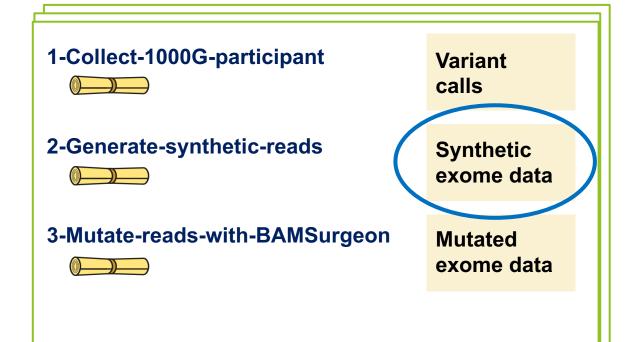Synthetic exome data

3-Mutate-reads-with-BAMSurgeon

Mutated exome data

How do we turn this into a FAIR community resource to empower biomedical researchers to leverage the underlying tools more easily?

# This hackathon project broad.io/FAIRdatahack2019

1-Collect-1000G-participant

2-Generate-synthetic-reads

3-Mutate-reads-with-BAMSurgeon

Variant calls

Synthetic exome data

Mutated exome data

**#1 - Data in demand**

What kind of datasets would be useful to the community?

Identified top needs based on literature + feedback from hackathon participants

# The Workspace

# The Next Steps

Off-the-shelf synthetic data catalog

User-friendly tools for generating custom synthetic (sequence) datasets

# The Team

## Broadies

Adelaide Rhodes
Allie Hajian
Anton Kovalsky
Ruchi Munshi
Tiffany Miller
Geraldine Van der Auwera

## Guest stars

Ernesto Andrianantoandro
Dan Rozelle
Jay Moore
Rory Davidson
Roma Kurilov
Vrinda Pareek