# PRECISION Medicine 2019
## Can AI Accelerate Precision Medicine?

## Analysis of scalable genomic variant data with HAIL using clusters in cloud environments

Instructors: Ines Krissaane, Lindsay Leek and Carlos De Niz

HARVARD
MEDICAL SCHOOL

# Contents

# Variants in the Human Genome

- The human genome is made of 3.2 billion $(3 \times 10^9)$ base pairs, distributed across 20k genes and comprised in 23 pairs of chromosomes

- Only $\sim$1.5% of the genome codes for protein

- DNA-wise, humans are 99.9% the same. The remaining 0.1% of these genetic variations, or **variants**, may give place to differences such as distinct skin, eye or hair color, but some of them may also be linked to or causative of disease



U.S. National Library of Medicine

- Human mutation rate is estimated to be $2.5 \times 10^{-8}$ mutations/bp/generation, which leads to 100-200 new variants in each newborn

- It is crucial to study and understand the possible effects of variants towards increasing risk for disease

# The Variant Call Format (VCF)

**Aim:** to score the similarity of alleles (SNPs, indels) to reference genome in scalable data sets ▸ Click HERE for additional information  ▸ VCF 4.1 Spec



Sections of a VCF file

# Scalability

- Genomic data grows fast. Analyzing it requires the use of clusters as often single machines can no longer process or store all these large data

- Clusters are flexible, as one can increase the number of working machines (cores) to easily boost computation power

- Data and system scalability involves virtual cluster deployment, monitoring and management in cloud environments

- An example of data scalability is the 1,000 Genome project, which has evolved ever since it was created:

| 1000 Genome Release | Variants | Individuals | Populations |
|---|---|---|---|
| Phase 3 | 84.4 million | 2504 | 26 |
| Phase 1 | 37.9 million | 1092 | 14 |
| Pilot 1 | 14.8 million | 179 | 4 |

# Hail, a genomic variant store

- Open-source, modular, scalable platform for statistical genetics in continuous development by the Neale lab at the Broad Institute
- It is exposed through Python and backed by distributed algorithms built on top of Apache Spark



Python API for Hail

# From VCF to MT file

A **Matrix Table (MT)** is a Hail specific data representation made by three tables: rows, columns and entries. Variant data is transformed into .mt format for processing in Hail. Sa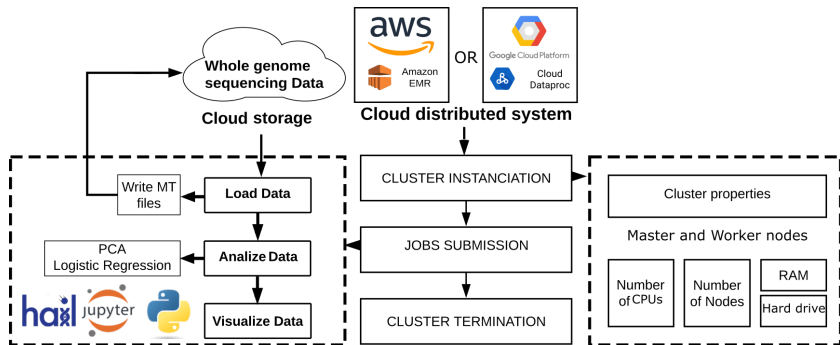mples are found in the columns table, alleles in the rows and format information from the VCF file in the entries table.



Hail Matrix Table Format

# A distributed computational framework for large scale genomic analysis.

# Additional Information

**Documentation and support:**

- Documentation, tutorials and sample code:
  hail.is
- Forum and chat:
  discuss.hail.is

**Spinning clusters in cloud services:**

- Hail Deployment in Google Cloud:
  github.com/hms-dbmi/Hail-on-Google-Cloud
- Hail Deployment in Amazon Web Services:
  github.com/hms-dbmi/hail-on-AWS-spot-instances