

# Analyzing Large Scale Genomic Data Using the Google Cloud Platform

**Cuiping Pan, Ph.D.**

Palo Alto Veterans Institute for Research  
VA Palo Alto; [Cuiping@stanford.edu](mailto:Cuiping@stanford.edu)

# 460 Human Genomes from MVP Pilot

## One Genome

- 3 billion DNA base pairs
- 3 million single nucleotide variation
- 50,000 private variants

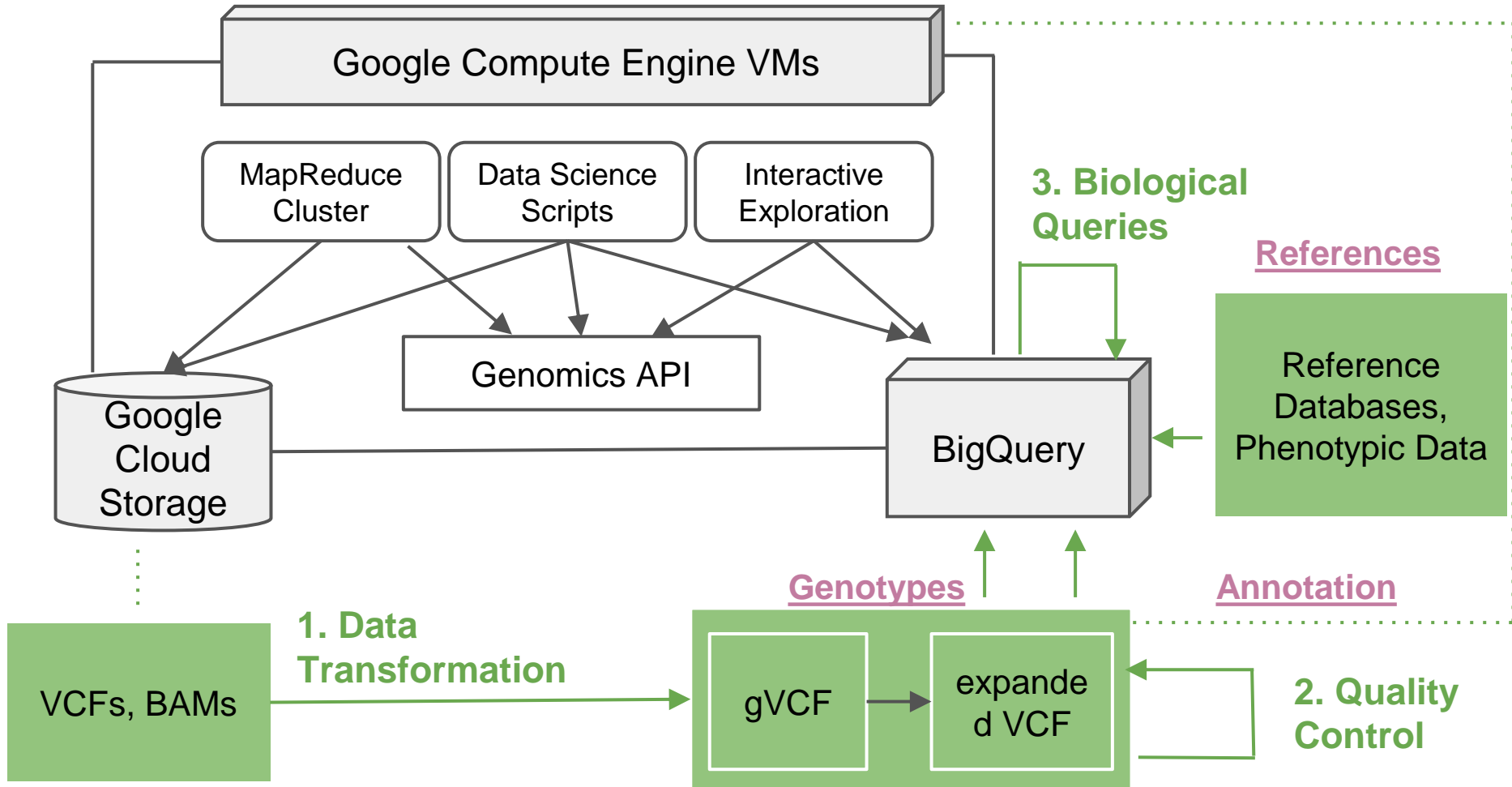
## 460 Genomes

- 48T pre-processed data
- 26M single nucleotide variation
- 5M short insertion and deletion

**Million Veteran Program:  
A Partnership with Veterans**



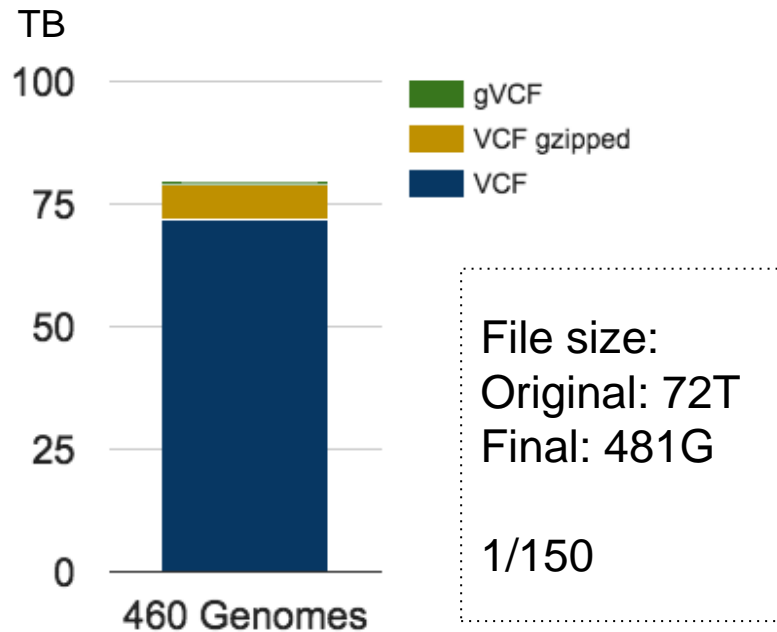
# GCP and Our Workflow



# Data Transformation and BigQuery Schema

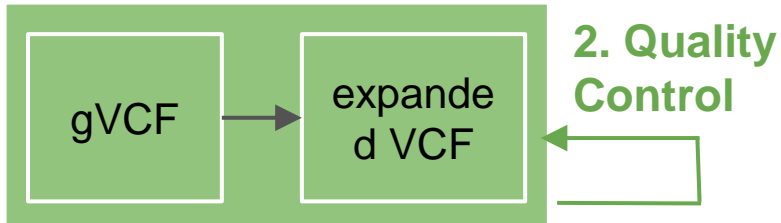
## gVCF:

- filter/tag variants by quality
- collapse consecutive REFs



reference_name	STRING
start	INTEGER
end	INTEGER
reference_bases	STRING
alternate_bases	STRING
call	RECORD
-- call.call_set_id	STRING
-- call.call_set_name	STRING
-- call.genotype	INTEGER
-- call.FILTER	STRING
-- call.QUAL	FLOAT
...	

# Multi-level Quality Controls



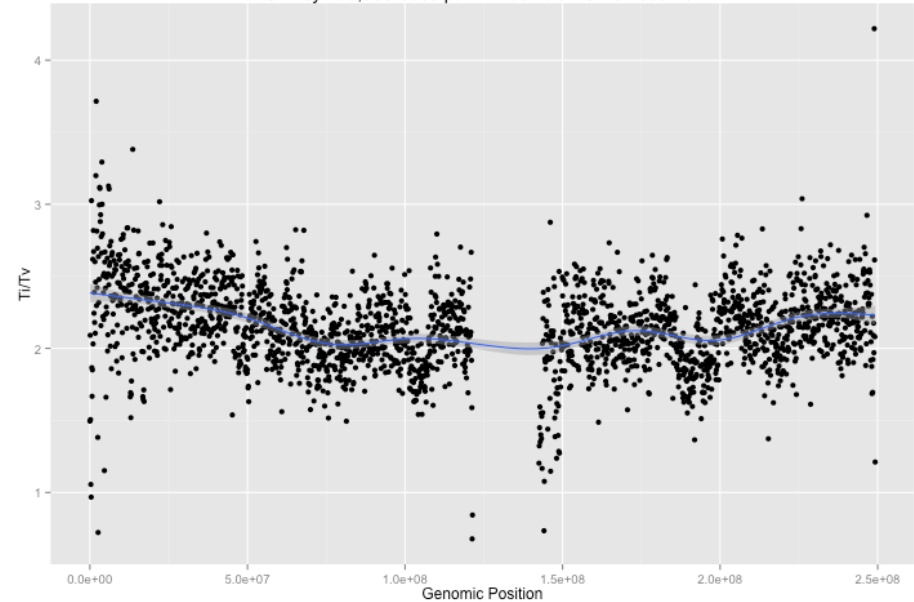
## Sample Level QC

1. batch effect
2. concordance to SNP array data
3. sex inference
4. ethnicity inference
5. genome similarity: IBS
6. missingness rate
7. singleton rate
8. heterozygosity rate
9. inbreeding coefficient

## Variant Level QC

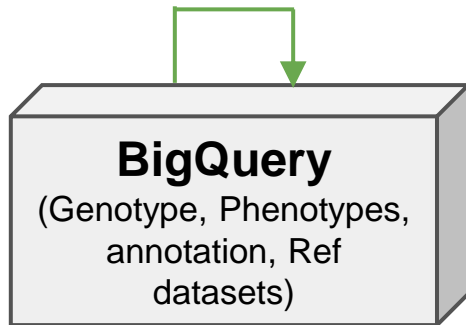
1. remove variants on blacklisted genes
2. heterozygous haplotype
3. missingness rate
4. Hardy-Weinberg Equilibrium
5. Ti/Tv by depth
6. Ti/Tv by alternate allele counts

Ti/Tv by 100,000 base pair windows on Chromosome 1

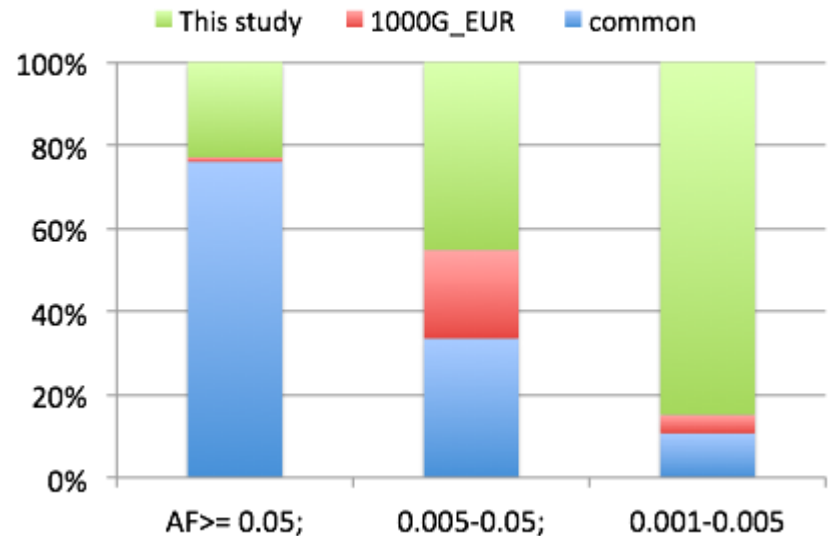


# Biological Queries

## 3. Biological Queries

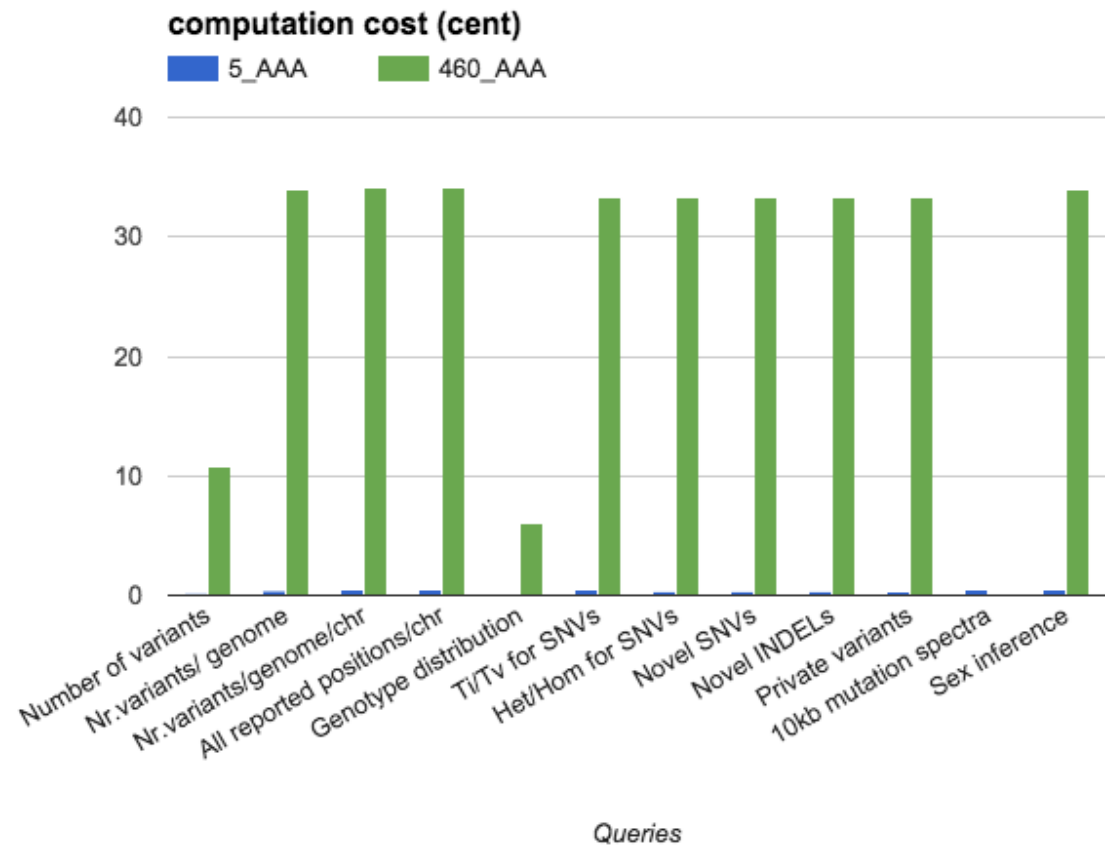
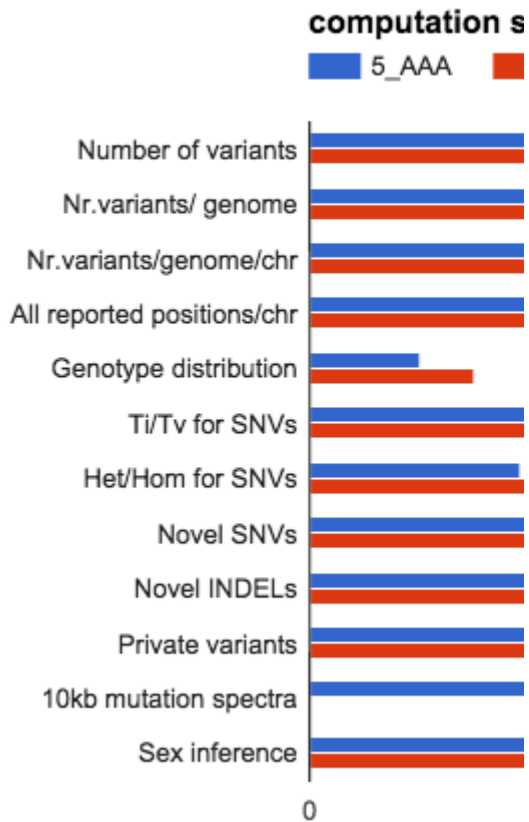


- Querying Genotypes
- Querying Biological Features
- Querying Medical Information



<i>VKORC1</i> Genotype (-1639G>A, <a href="#">rs9923231</a> )	<i>CYP2C9</i> *1/*1	<i>CYP2C9</i> *1/*2	<i>CYP2C9</i> *1/*3	<i>CYP2C9</i> *2/*2	<i>CYP2C9</i> *2/*3	<i>CYP2C9</i> *3/*3
GG	5-7	5-7	3-4	3-4	3-4	0.5-2
GA	5-7	3-4	3-4	3-4	0.5-2	0.5-2
AA	3-4	3-4	0.5-2	0.5-2	0.5-2	0.5-2

# Scalability and Affordability of GCP



# Acknowledgements

## VA Palo Alto

Philip Tsao  
Alicia Deng

## Stanford University

Gregory McInnes  
Somalee Datta  
Isaac Liao  
Denis Salins  
Michael Snyder

## Google Genomics

Nicole Deflaux  
Jonathan Bingham  
Elmer Garduno  
David Glazer  
Samuel Gross

Poster #1w

Email: [cuiping@stanford.edu](mailto:cuiping@stanford.edu)