

# TruCal Variant Recalibrator

---

SVAI Hackathon 2017

Team: Mutrakers

Problem Track: ranking dataset mutations

# Team

**Don Freed, Ph.D**  
Bioinformatics Scientist



**Annabelle Tang**  
Research Associate



**Alex Francis**  
Data Scientist



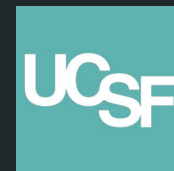
**David Streid**  
Software Engineer



**Magdalena Matusiak, Ph.D**  
Postdoc Fellow Bioinformatics

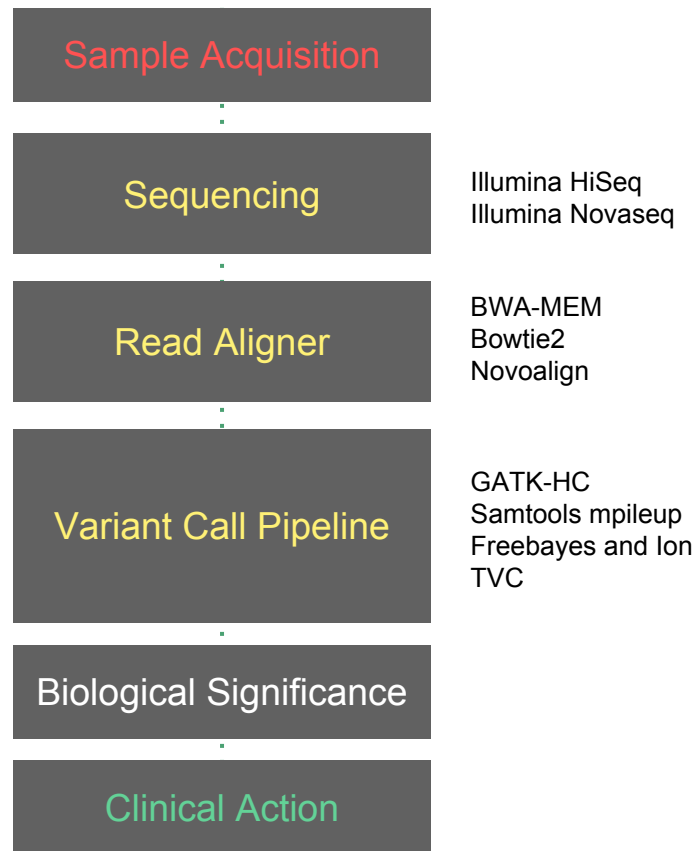


**Chris Margono**  
Medical Student



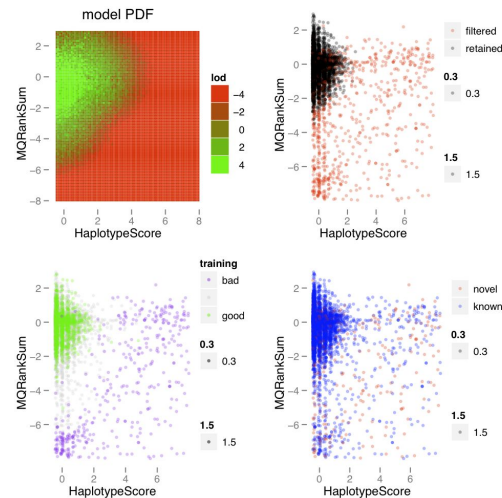
# Introduction

- Translating data to care requires:
  - proper processing
  - Proper interpretation
- **Errors** pose a risk to clinical application
- It has been demonstrated that different variant call methods can produce discordance outputs, indicating the need for careful interpretation of their results

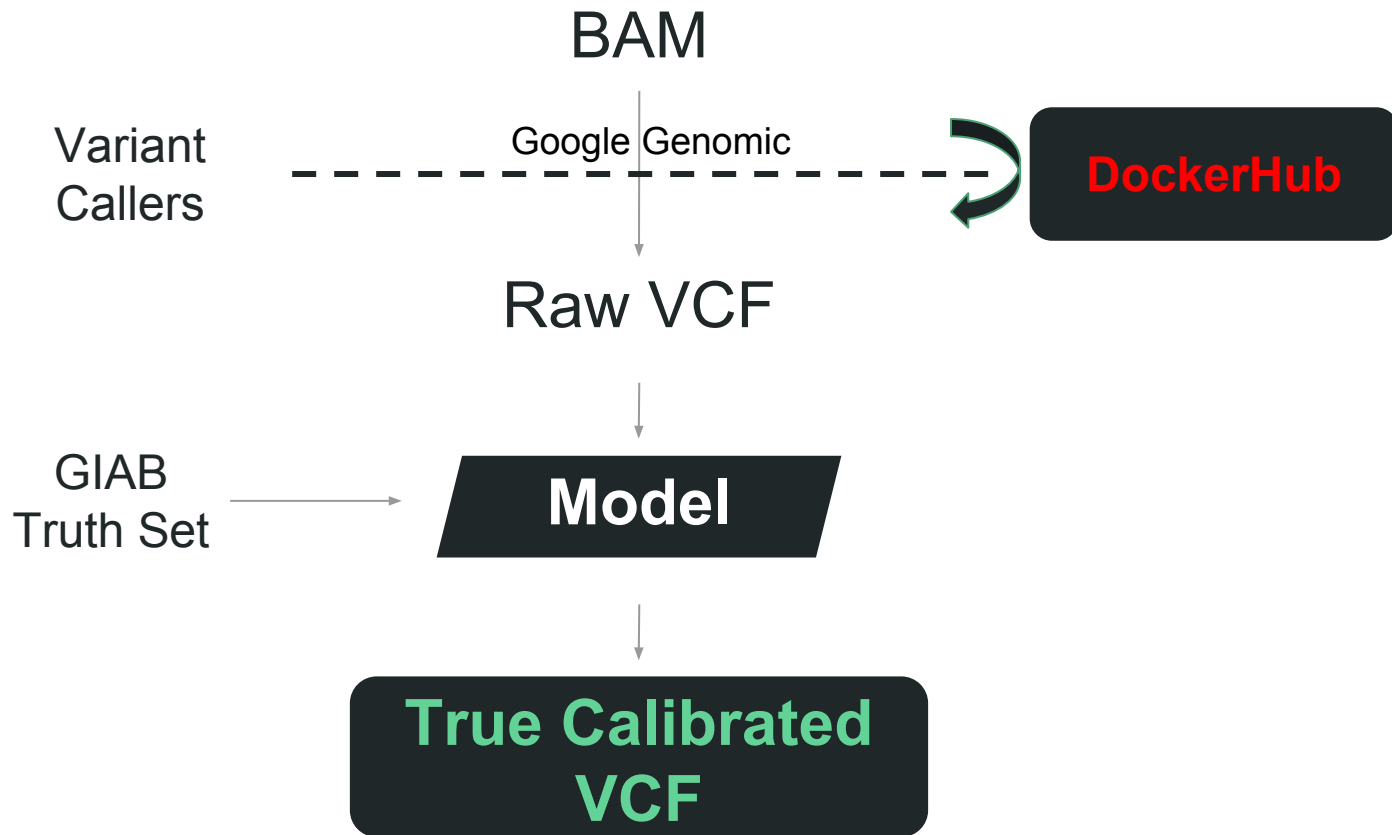


# Purpose

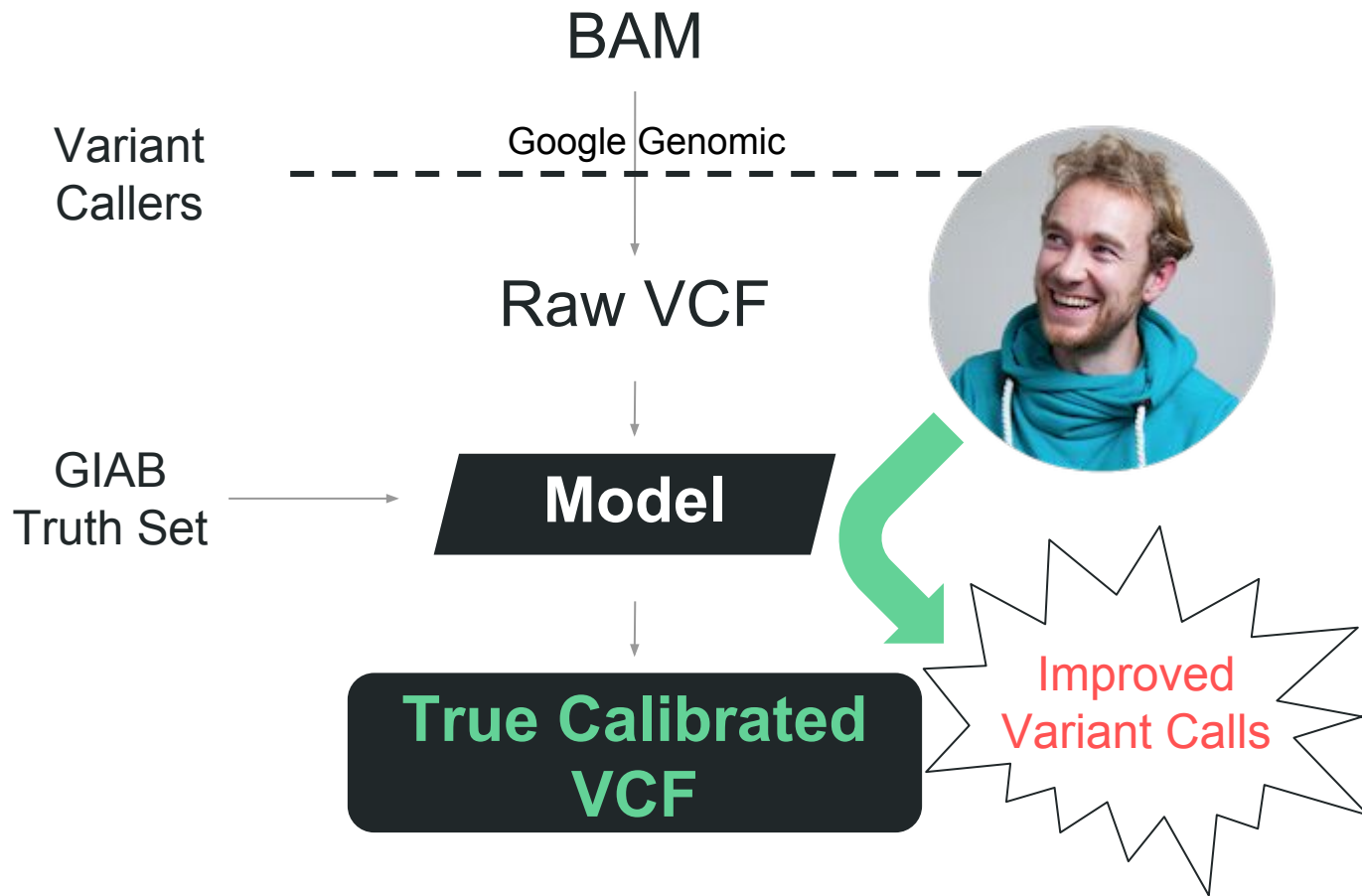
- Previous work
  - used unsupervised learning to sort true mutations and errors from variant call outputs
  - no gold-standard truth set existed at that time
- Recent advancements (2015)
  - Genome in a Bottle Consortium (GIAB) Released high-confidence variant calls
- **Goal:**
  - **improve precision** (positive predictive value) of variant call outputs using **supervised learning** with **GIAB data**



# Approach



# Impact



# Methods: Data Transformations

$$X = \begin{matrix} & \overbrace{\hspace{1cm}}^{\text{AD}} & & \overbrace{\hspace{1cm}}^{\text{QUAL}} \\ \begin{bmatrix} x_{11} & \dots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{1d} \end{bmatrix} \end{matrix}$$

*From the “True VCF,” utilize a “featurizer” to obtain a data matrix for training*

- *Impute missing features (messy data!) using the median of all features.*
- *Labels (“real variants” boolean) are compiled from Genome in a Bottle project*

$n \approx 10^6$  (the number of reads/“mutations” in the VCF),  $d = 21$

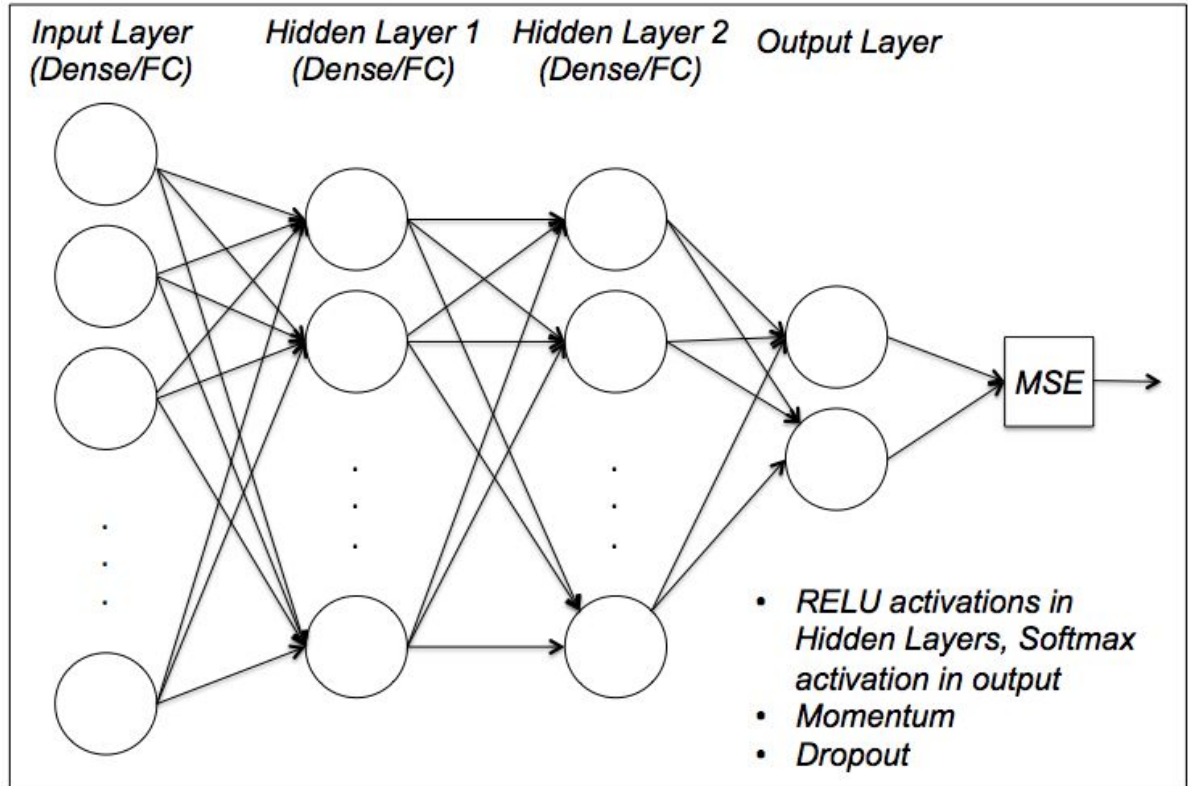
features = {

AD, GQ, DP, GT, ExcessHet, AC, BaseQRankSum, FS, AF, MLEAC, AN, SOR,  
MQ, QD, DP, ClippingRankSum, MQRankSum, ReadPosRankSum, QUAL, is\_snp

}

# Methods: Machine Learning

- Train three-layer neural network using Keras
- Produces a “probability calculator” from the output of the softmax activation in the final layer





# Results

- Check out our model!

<https://github.com/SVAI/MutRackers>