

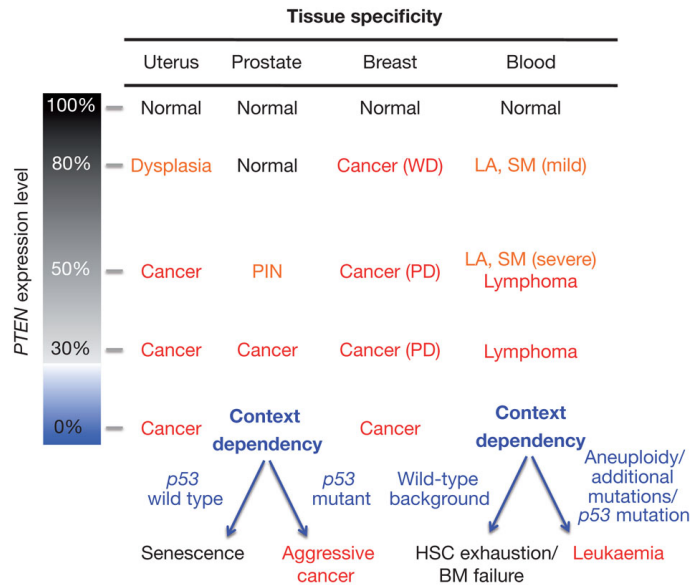
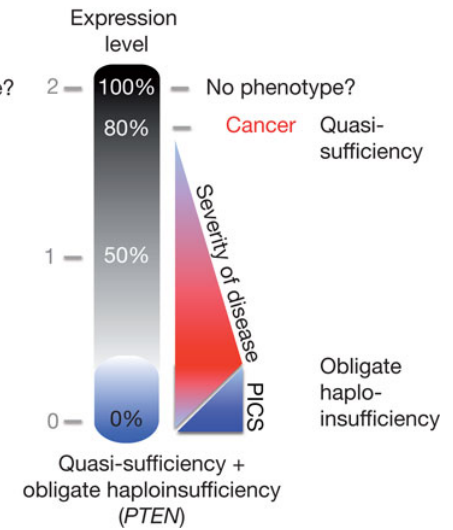
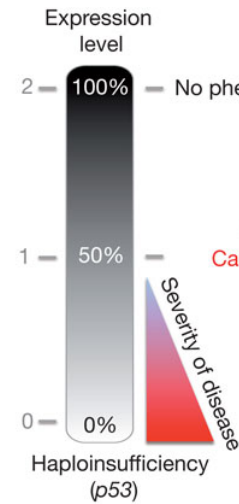
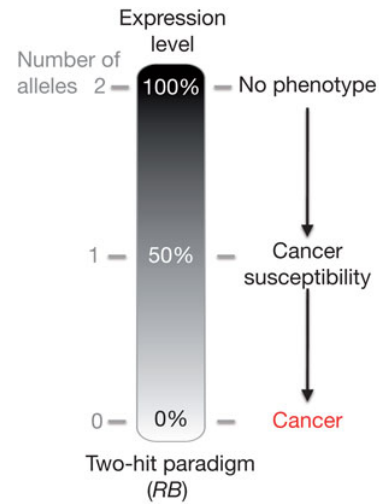
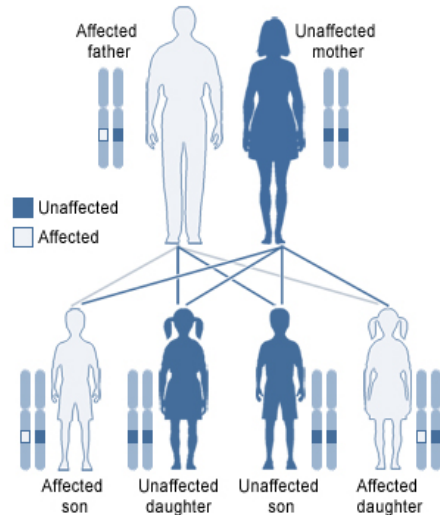


# Improving Next-Generation Sequencing Variants Identification In Cancer Genes Using Globus Genomics

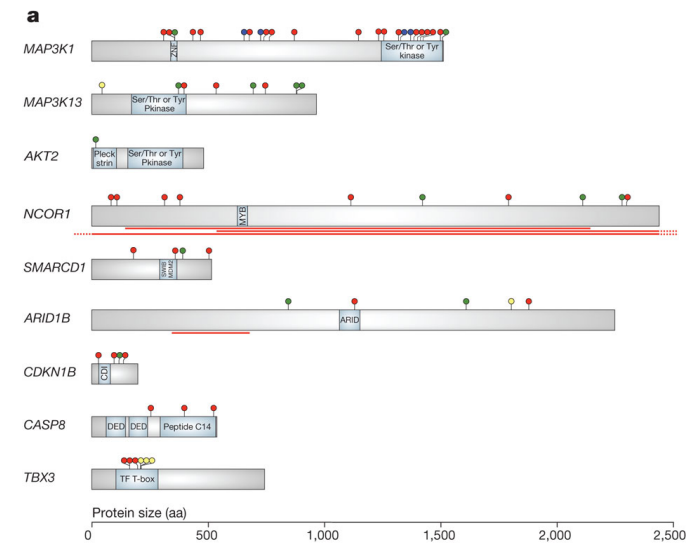
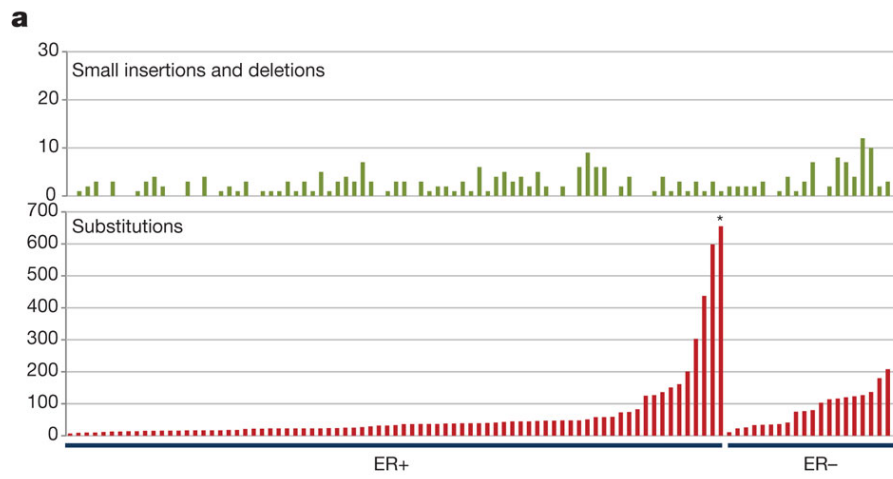
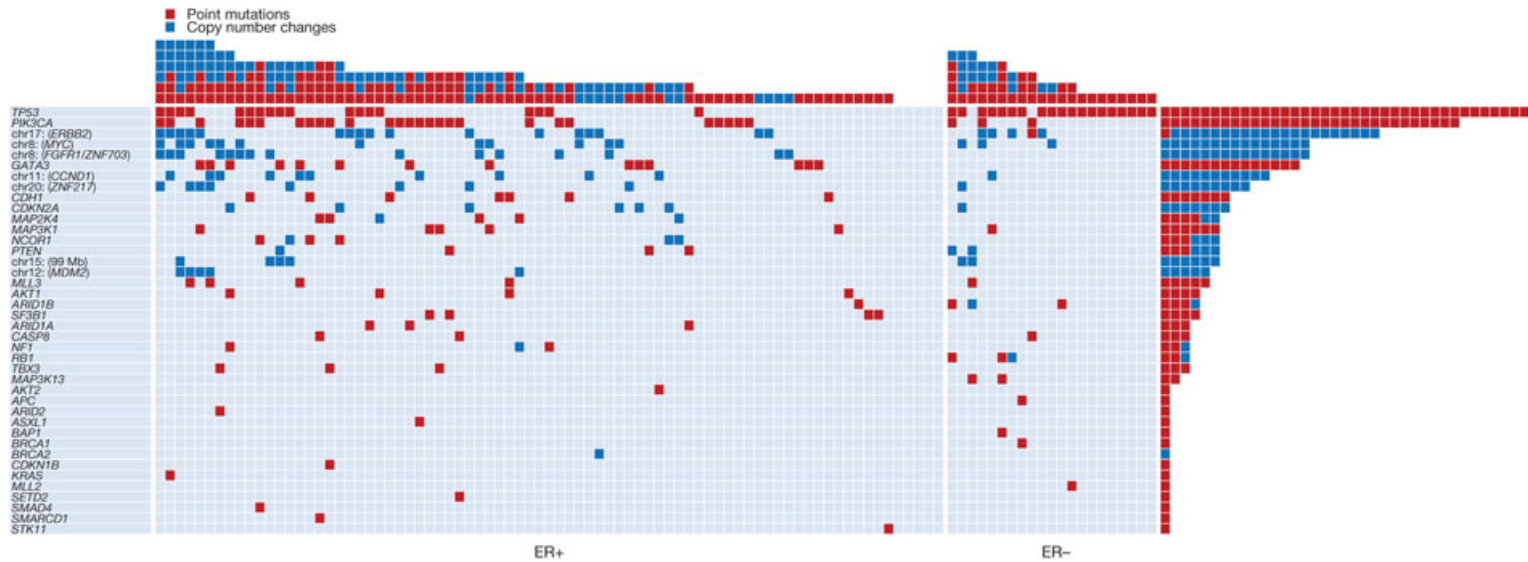
Toshio F Yoshimatsu<sup>1</sup>, Yonglan Zheng<sup>1</sup>, Alex Rodriguez<sup>3</sup>, Vassily Trubetskoy<sup>2</sup>, Ravi  
K Madduri<sup>3</sup>, Paul J Dave<sup>3</sup>, Nancy J Cox<sup>2</sup>, Ian T Foster<sup>3</sup>, Olufunmilayo I Olopade<sup>1</sup>

1. Department of Medicine, Section of Hematology/Oncology
  2. Department of Medicine, Section of Genetic Medicine
  3. Computation Institute
- The University of Chicago

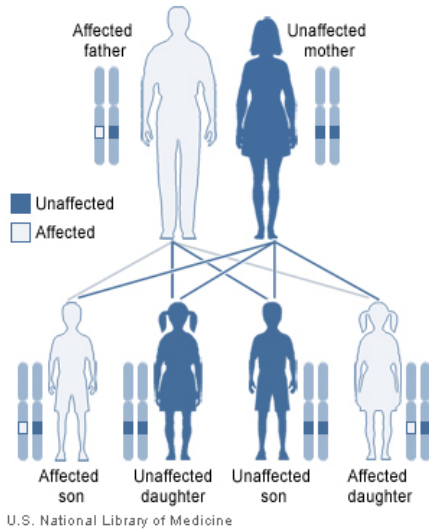
# Cancer is a genetic disease



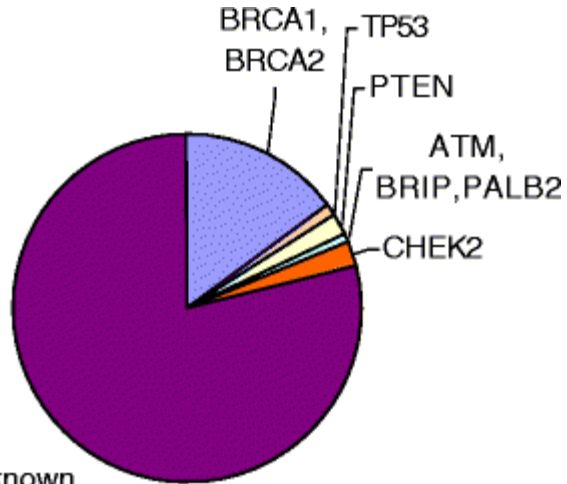
# Breast Cancer is extremely heterogeneous



# Strategies for Breast Cancer risk assessment

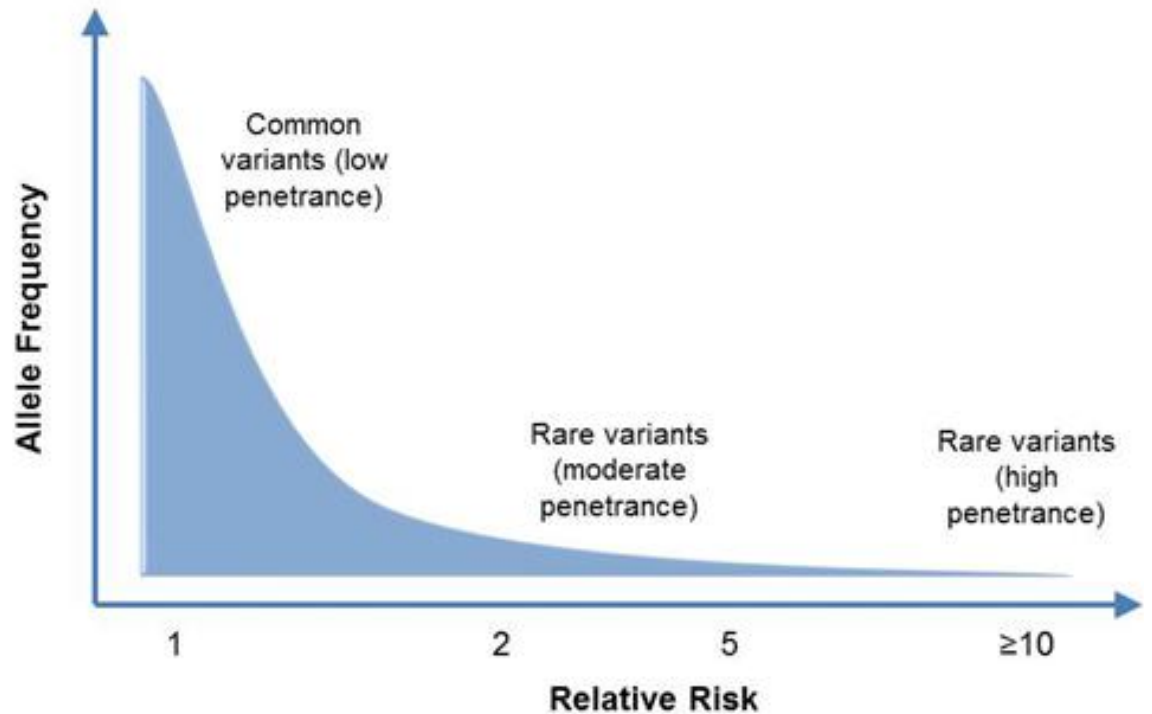


U.S. National Library of Medicine



Unknown

## Genetic Architecture of Cancer Risk



Van der Groep *et al. Cell Oncol* **34**, 71-88 (2011).

# BROCA Cancer Risk Panel

Developed by Dr. Mary-Claire King and colleagues at the University of Washington at Seattle

“This assay sequences all exons and flanking intronic sequences of 50 cancer genes. A total of 1.1 Mb (1.1 Million base pairs) are sequenced and the average coverage ranges from 320 to >1,000 sequencing reads per bp. Genomic regions are captured using biotinylated RNA oligonucleotides (SureSelect), prepared in paired-end libraries with ~200 bp insert size, and sequenced on an Illumina HiSeq2000 instrument with 100 bp read lengths.”

(<http://web.labmed.washington.edu/tests/genetics/BROCA>)

# BROCA in 200 Nigerian breast cancer cases

<u>Phenotype</u>	<u>Number</u>
early age onset (<50yrs)	158
familial	20
familial & early age onset	19
others	3

# Real mutations are hidden in the noise



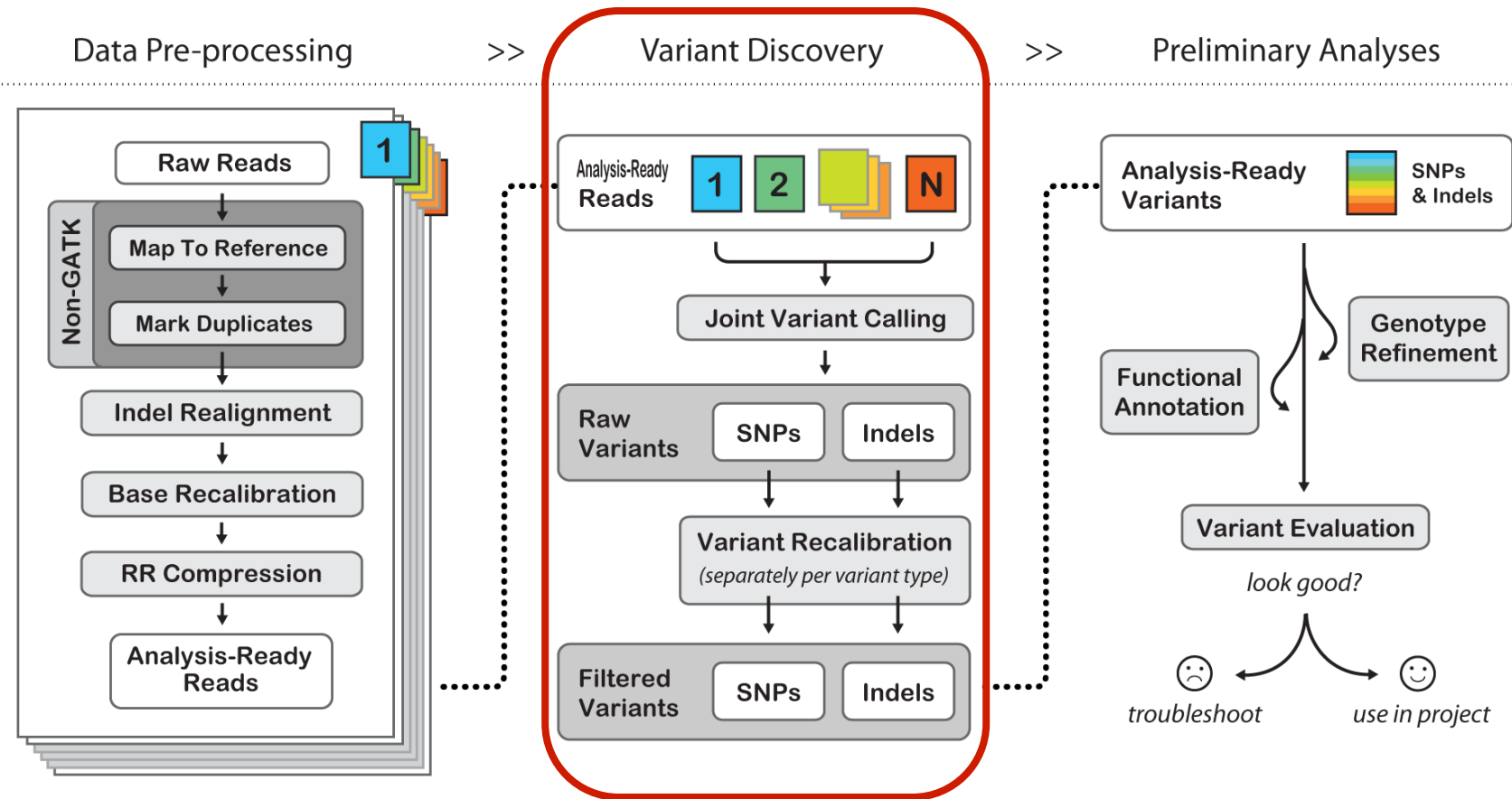
How to tell which mismatch is real mutation and which are just noise?  
i.e. How can we reduce false positive and false negative detection rates?

# Motivation

We want to improve sensitivity and specificity of variant identification using next-generation sequencing data

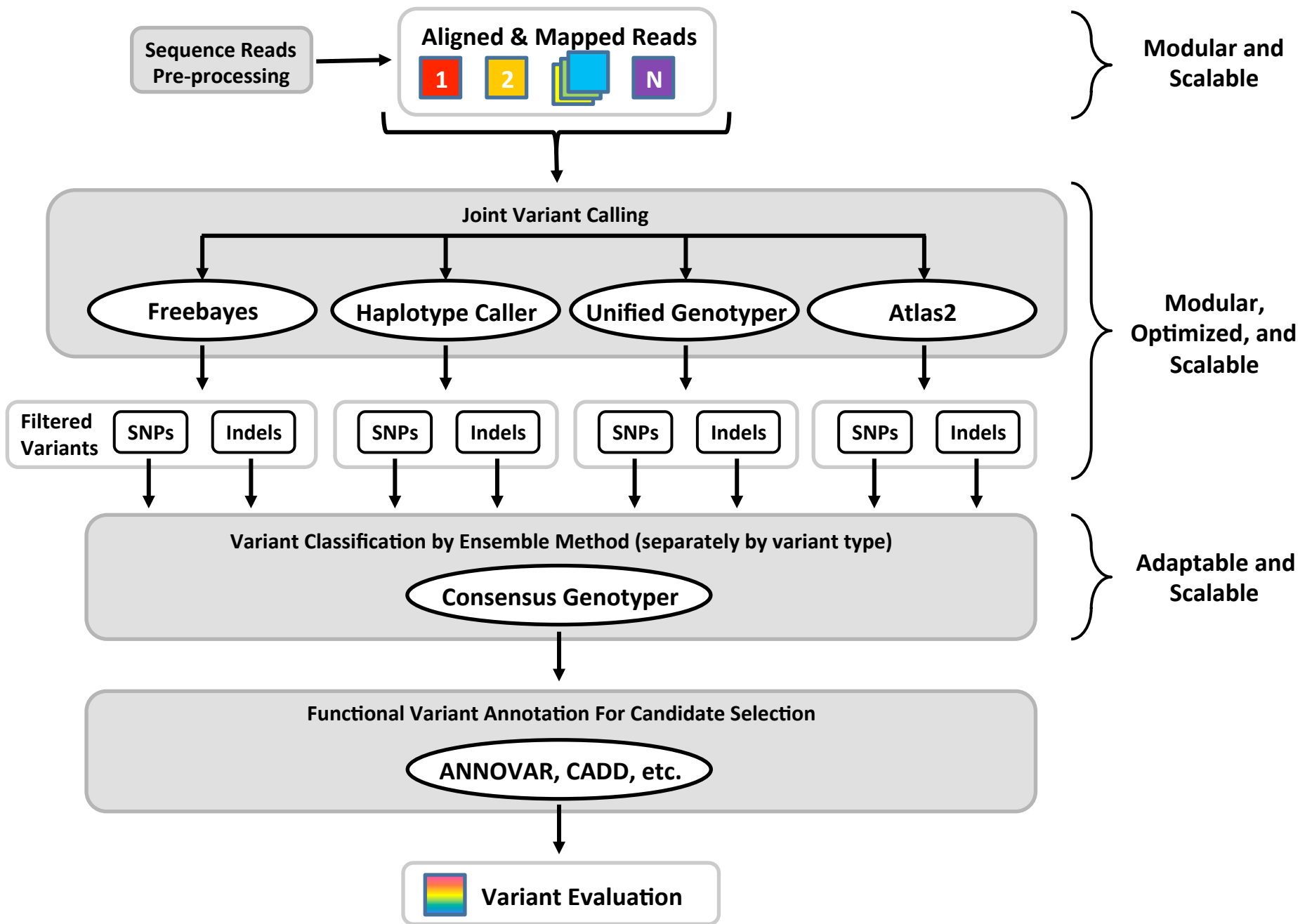


# GATK Best Practice



Our work focuses on this step

\*Image courtesy of Broad Institute GATK Team  
(<http://www.broadinstitute.org/gatk/guide/best-practices>)



# Computation Performance

**TABLE 2.** Summary for the alignment of 200 BROCA target exome-seq Fastq files.

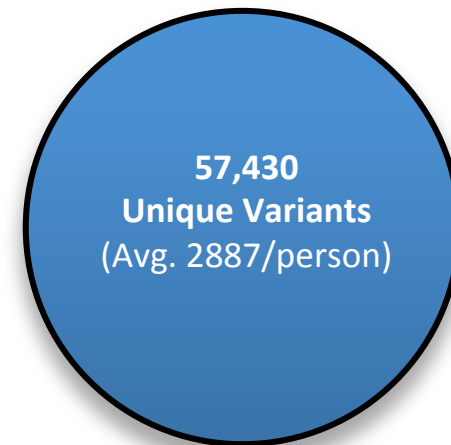
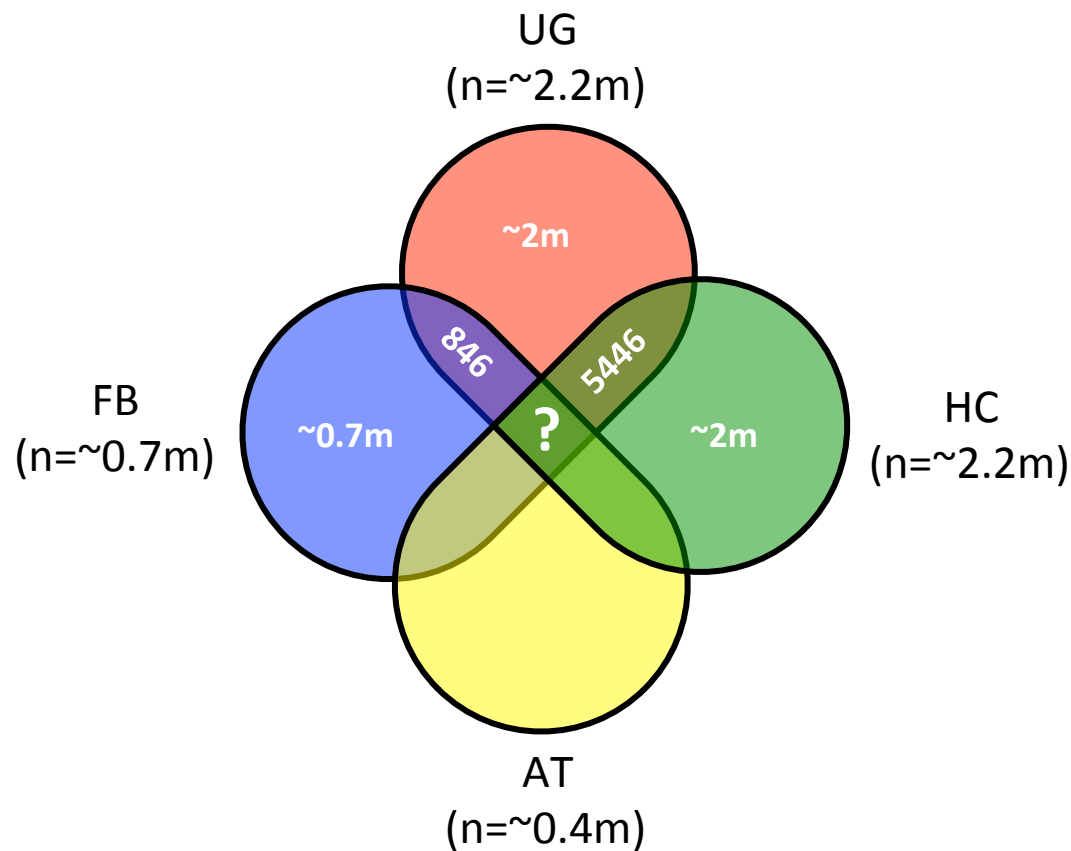
	<b>GATK UG</b>	<b>GATK HC</b>	<b>Freebayes</b>	<b>ATLAS2-SNP</b>
<b>Input Type</b>	BAM	BAM	BAM	BAM
<b><u>Paralellization Level</u></b>	Chromosome (24X)	Chromosome (24X)	Chromosome (24X)	Input File (200X)
<b>Input Size</b>	67.5 GB	67.5 GB	67.5 GB	67.5 GB
<b>Data Generated</b>	54 GB	19 GB	0.70 GB	25 GB
<b>Output Size</b>	615 MB	615 MB	315 MB	1.5 GB
<b>Total CPU time</b>	22 hours	27.3 hours	37.2 hours	708 hours
<b>Walltime for analysis</b>	6.75 hours	2.33 hours	4.33 hours	23 hours
<b>Worker Nodes Used</b>	13	23	24	200

\*Optimizations can still be achieved by running multiple chromosomes and samples on the same worker node

# Variant Calling Performance - SNPs

## Globus Genomics

## BROCA Default

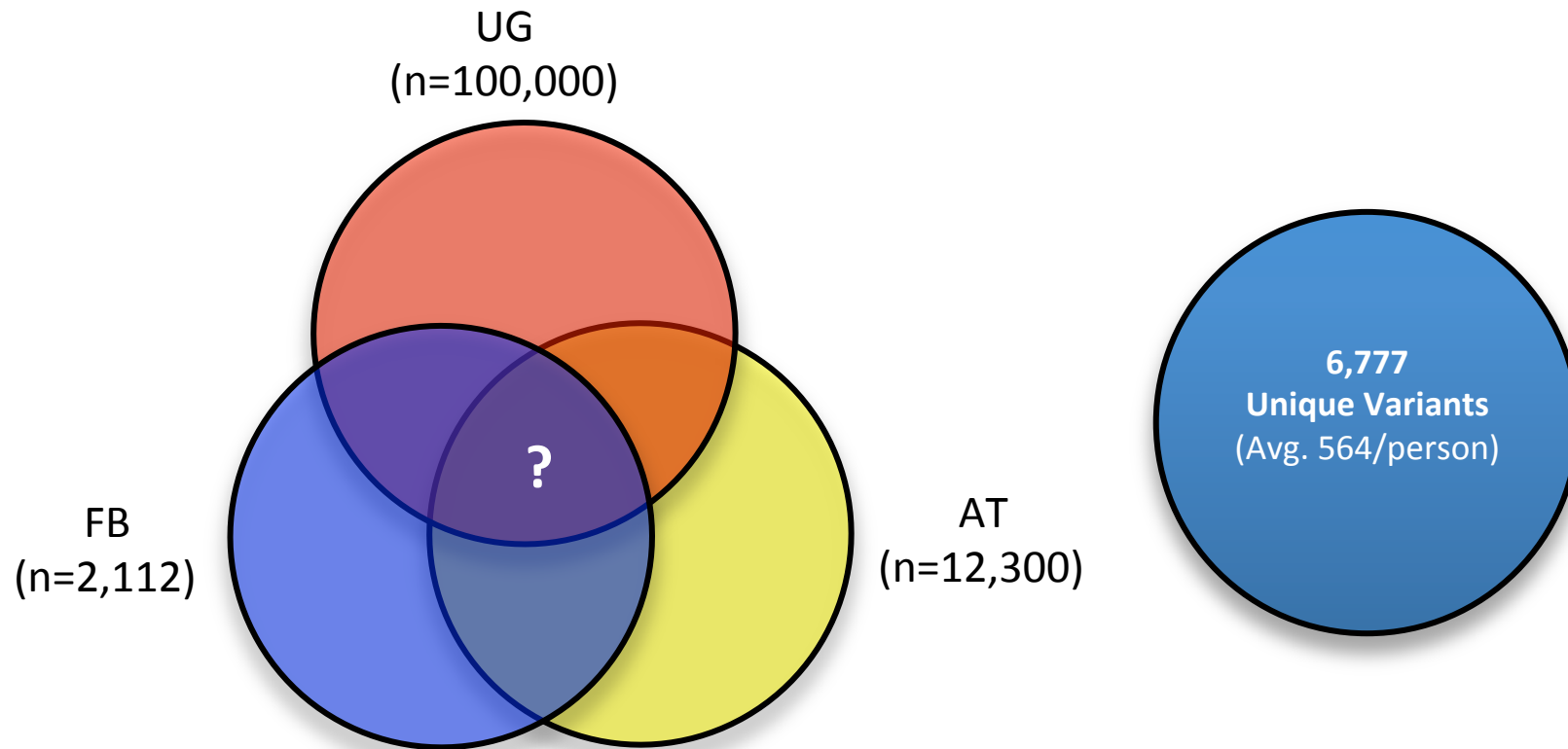


UG = Unified Genotyper; HC = Haplotype Caller; FB = Freebayes; AT = Atlas2

# Variant Calling Performance - Indels

## Globus Genomics

## BROCA Default



UG = Unified Genotyper; FB = Freebayes; AT = Atlas2

# Summary

- We believe our variant calling approach can be easily adapted and modified for generalization and implementation in other genome-wide and large-scale NGS variant analysis.
- Our method fully utilizes the capacity of Globus Genomics to make the workflow scalable, modular, and adaptable. Consequently, analysis time were dramatically shortened (often in the order of magnitude).
- Incorporation of Globus Online into the pipeline automates and facilitates transfer and sharing of large data.
- We are further developing the previous analysis method to include other variation types into consideration (e.g. indels, copy number change).
- Our ultimate goal is to help patients make informed decision, and not to scare them with false alarm.

# Acknowledgements



## University of Chicago

*Dept. Medicine, Hem/Onc*

Niu Qun

**Yonglan Zheng**

Funmi Olopade

*Dept. Medicine, Genetic Medicine*

Vassily Trubetskoy

Nancy Cox

*Computation Institute*

**Alex Rodriguez**

Ravi Madduri

Paul Dave

Ian Foster



## University of Washington at Seattle

Mary-Claire King

Tom Walsh

Ming Lee

## University of Ibadan

Oladosu Ojengbede

Temidayo Ogundiran

Abideen Oluwasola

Abayomi Odetunde

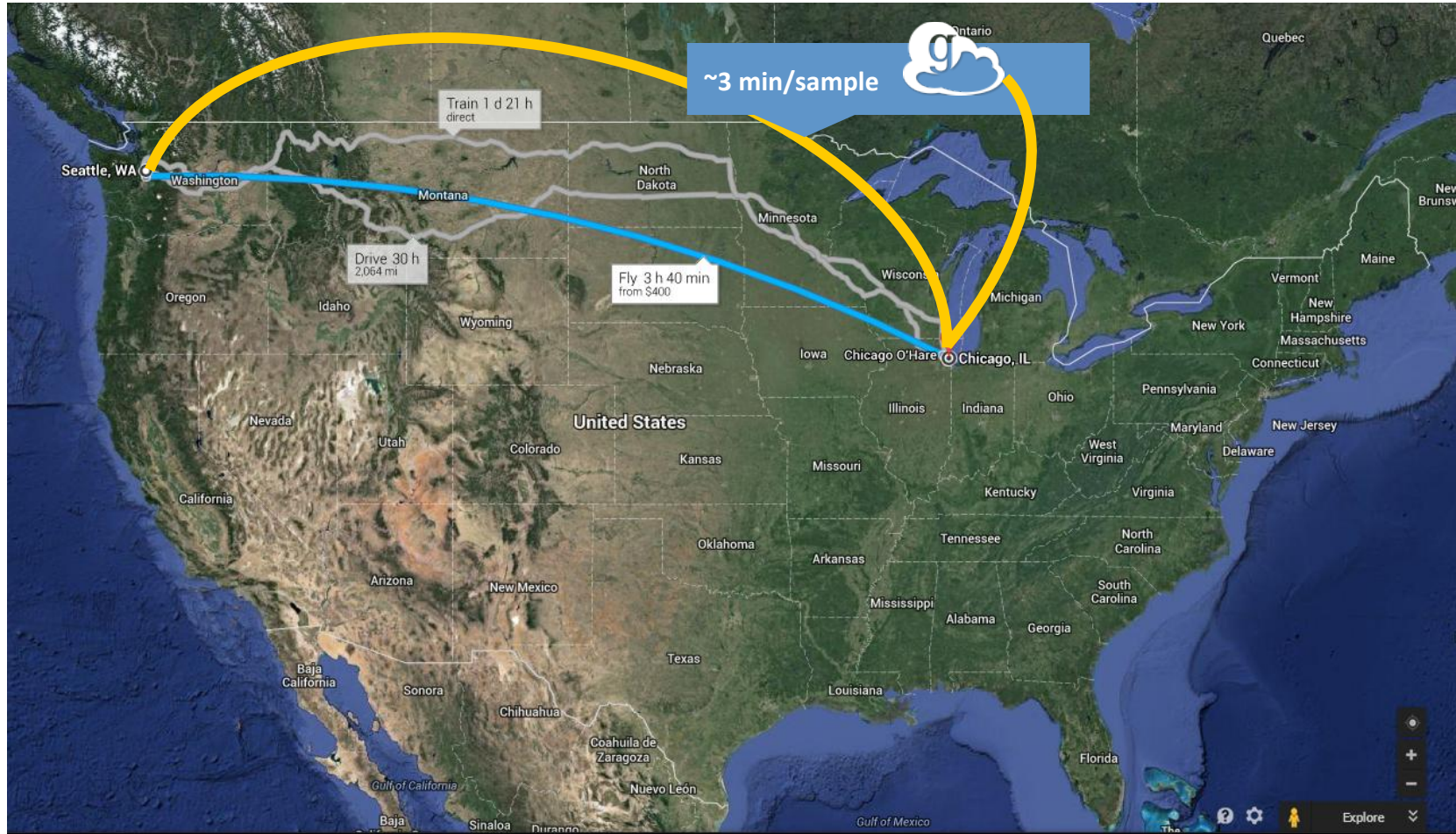
## Support

Susan G Komen for the Cure

Breast Cancer Research Foundation

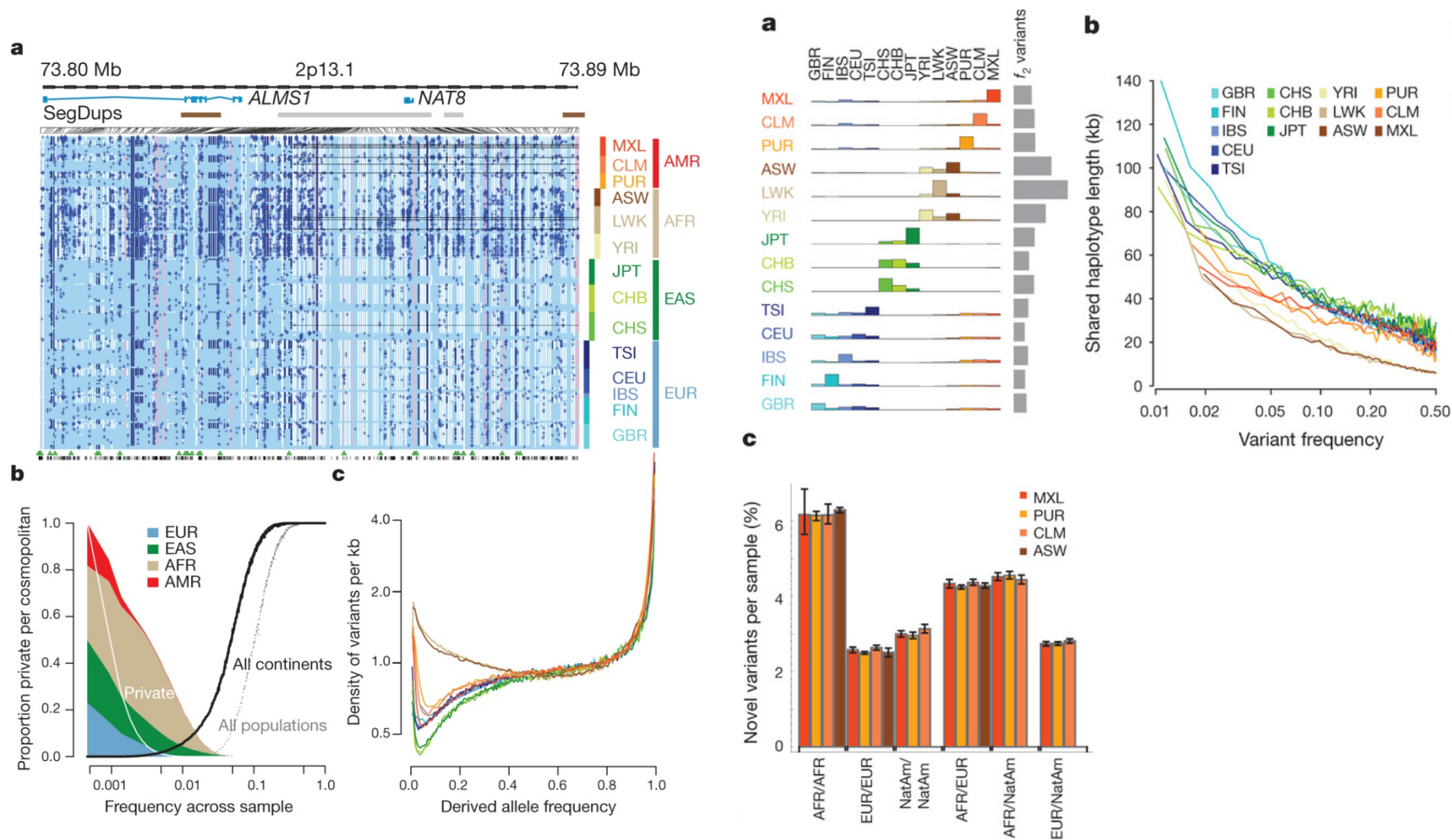


# Data Transfer

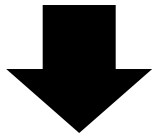




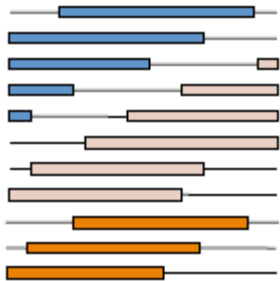
# Africans are known to have complex genetic background



# General workflow of NGS



Sequencing  
Output



Enormous pile  
of short reads  
from NGS



Cleaning and  
Organizing Data

We are searching for variants

