

Overview of the NHLBI Trans-Omics for Precision Medicine (TOPMed) program: whole genome sequencing of >100,000 deeply phenotyped individuals

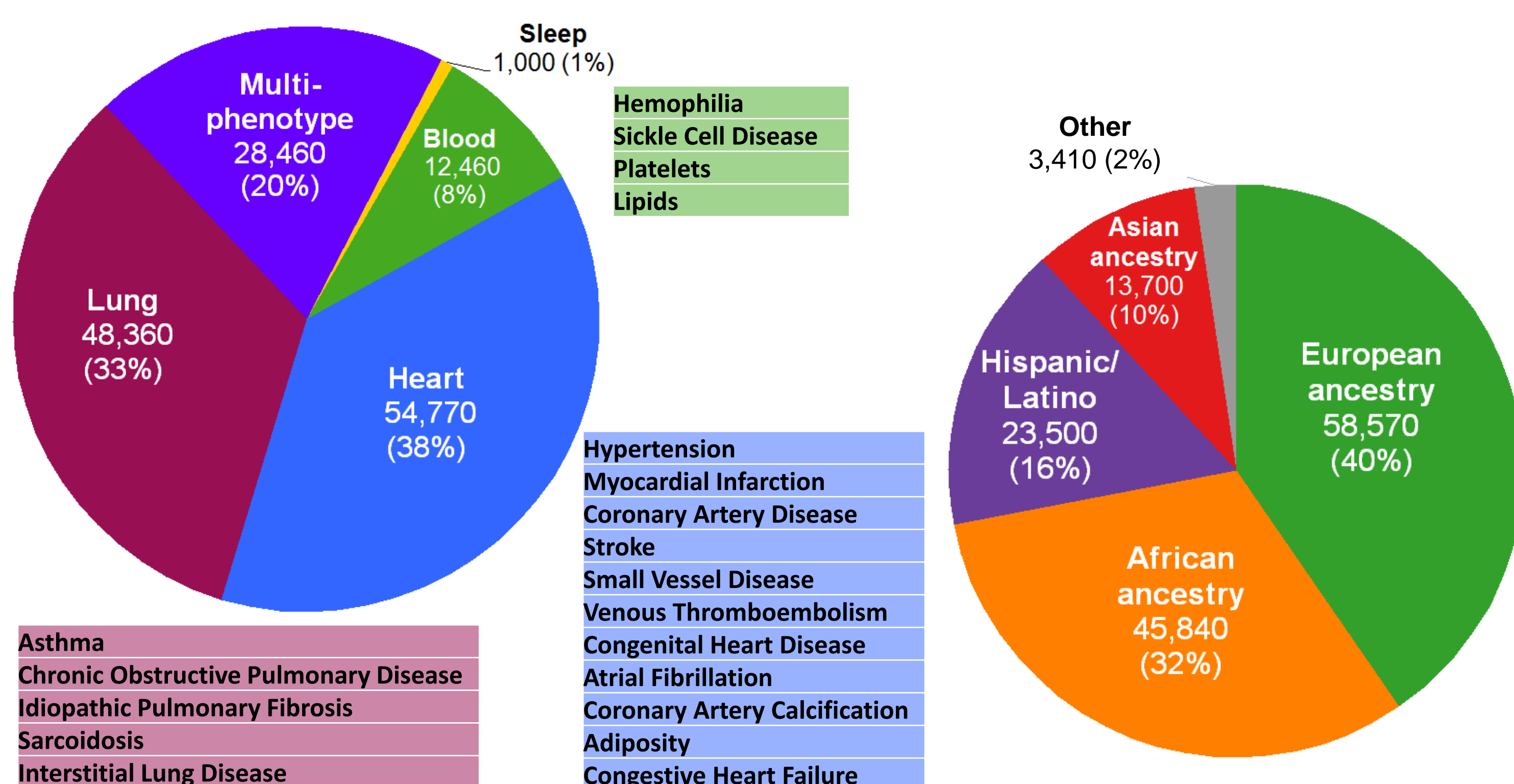
Cathy Laurie¹, Tom Blackwell², Goncalo Abecasis², Ken Rice¹, James Wilson³, Deborah Nickerson⁴, Stacey Gabriel⁵, Richard Gibbs⁶, Susan Dutcher⁷, Soren Germer⁸, Donna Arnett⁹, Allison Ashley-Koch¹⁰, Kathleen Barnes¹¹, Eric Boerwinkle¹², Stephen Rich¹³, Edwin Silverman¹⁴, Rebecca Beer¹⁵, Julie Mikulla¹⁵, Pothur Srinivas¹⁵, Weiniu Gan¹⁵, George Papanicolaou¹⁵, Cashell Jaquish¹⁵, and the NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium

1) Department of Biostatistics, University of Washington, 2) Department of Biostatistics, University of Michigan, 3) Department of Physiology, University of Mississippi, 4) Department of Genome Sciences, University of Washington, 5) Broad Institute, Massachusetts Institute of Technology, 6) Baylor College of Medicine Human Genome Sequencing Center, 7) McDonnell Genome Institute, Washington University in St. Louis, 8) New York Genome Center, 9) Department of Epidemiology, University of Kentucky, 10) Department of Medicine, Duke University Medical Center, 11) Department of Medicine, University of Colorado at Denver, 12) Department of Epidemiology, University of Texas Health at Houston, 13) Center for Public Health Genomics, University of Virginia, 14) Channing Division of Network Medicine, Brigham & Women's Hospital, 15) National Heart, Lung and Blood Institute, NIH

Study Characteristics

A primary goal of the NHLBI TOPMed program is to improve scientific understanding of the fundamental biological processes that underlie heart, lung, blood, and sleep (HLBS) disorders. TOPMed is providing deep whole genome sequencing (WGS) and other omics data to pre-existing 'parent' studies having large samples of human subjects with rich phenotypic characterization and environmental exposure data.

Sample numbers by phenotype area and ancestry/ethnicity (N=144k total)

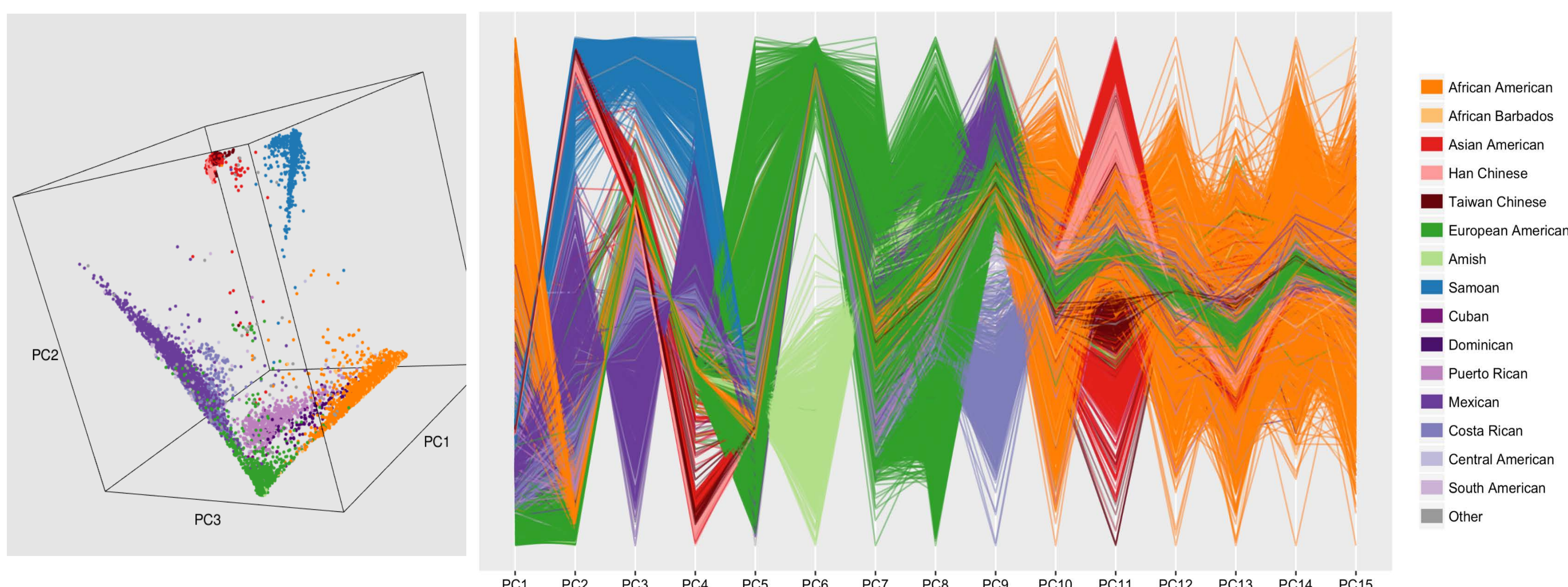


Study Designs

Currently, TOPMed consists of ~144k participants from >80 different studies with varying designs. Prospective cohorts provide large numbers of disease risk factors, subclinical disease measures, and incident disease cases; case-control studies provide large numbers of prevalent disease cases; extended family structures and population isolates provide improved power to detect rare variant effects. The phenotype pie chart above shows the numbers and percentages of participants in studies with a focus on HLBS, as well as the percentage belonging to cohort studies that have collected many different phenotypes.

Participant Diversity

Achieving ancestral and ethnic diversity was a priority in selecting contributing studies. Currently, the 144k participants consist of approximately **60% with substantial non-European ancestry** (see pie chart above, based on participant self-identification and study inclusion criteria). This diversity is also illustrated by the principal components analysis of ~55k participants (see 3D PCA plot and parallel coordinates plot below). Discovery of genotype-phenotype associations frequently includes pooled analysis across ancestry groups and studies, using statistical models that account for population structure and relatedness.



Acknowledgments

Whole genome sequencing (WGS) for the Trans-Omics in Precision Medicine (TOPMed) program was supported by the National Heart, Lung and Blood Institute (NHLBI). Please see <https://www.nhlbiwgs.org/acknowledgements> for full acknowledgments and funding sources. We acknowledge the support of Daniel Taliun, Quenna Wong, Sarah Nelson, Stephanie Gogarten, and Deepti Jain for providing data summaries and graphics in this poster.

Whole Genome Sequencing

Coverage, sequencing depth and number of variants in TOPMed freeze 5 (~65k samples)

	Total Number	% Singletons	Per Sample 5 th percentile	Per Sample median	Per Sample 95 th percentile
Bases (Gb)	8,523,992		107	131	156
Depth (x)			32	39	47
Genome covered %					
All			98.12	99.85	99.90
Depth > 10x			97.80	99.18	99.62
Total variants	471,006,757	46	3,514,025	3,567,972	4,363,040
SNVs	437,606,741	46	3,333,422	3,384,184	4,128,810
Indels	33,400,016	47	180,486	183,885	234,188
Coding variation					
Synonymous	1,734,456	43	10,837	11,078	13,688
Non-synonymous	3,644,788	48	10,632	10,880	13,231
Stop/essential splice	122,804	54	427	458	567
Indels					
Frameshift	134,489	60	111	128	167
In-frame	64,356	48	85	100	129

WGS was performed by several sequencing centers to a median depth of 39X using DNA from blood, PCR-free library construction and Illumina HiSeq X technology. A Support Vector Machine quality filter was trained with known variants and Mendelian-inconsistent variants. Genotypes were called jointly across all samples available to produce genotype data "freezes." For freeze 5, the median pair-wise non-reference discordance rates are: 4×10^{-4} for SNVs and 6×10^{-3} for indels passing the quality filter (N=378 duplicate sample pairs).

Explore TOPMed variants in Bravo

Bravo Variant Server Interface:

- Search: <https://bravo.sph.umich.edu/>
- Variant: rs2814778
- Gene: PCSK9, Transcript: ENST00000407236, Variant: chr22:16389447-A-G or rs34747326, Region: chr1:55030529-55075873
- Powered by Freeze5 on GRCh38
- The dataset includes 463 million variants on 62784 individuals
- Summary: Filter Status: PASS, Existing variation: rs2814778, Allele count: 32654/125568, Homozygous Alt Count: 12815, CADD: 16.8
- Frequency table: Population vs Allele Frequency (e.g., 1000G African: 0.9637)
- Site Quality Metrics: QC metric, Value, Percentile (e.g., TD: 2273030, 14.99)
- Annotation: This variant falls on 8 transcripts belonging to 3 genes (5'UTR, Intron, Downstream gene)
- IGV plots: Heterozygous Homozygous

TOPMed Resources for the Community

TOPMed data are being made available to the scientific community as a series of "data freezes": **genotypes and phenotypes via dbGaP**; read alignments via the SRA; and variants via the Bravo variant server and dbSNP. Genotypes for a set of 18.5k samples have been released on dbGaP, another freeze of 55k samples is currently being released, and a freeze of >100k samples is planned starting early 2019. TOPMed WGS data are contained in study-specific accessions with names containing "NHLBI TOPMed", while most phenotypic data are in parent study accessions. **The TOPMed web site has much more information about the study, as well as how to access the data through dbGaP — see www.nhlbiwgs.org.**

Other Omics

TOPMed is currently adding other omic assays to samples that have been whole genome sequenced; these include RNAseq, metabolomics, proteomics and epigenomics. These data will become available via dbGaP starting in 2019.