COVERING ALL THE BASES: A PRIMER ON TODAY'S SEQUENCING TECHNOLOGIES AND THEIR APPLICATIONS IN PRECISION MEDICINE RESEARCH

INSTITUTE FOR COMPUTATIONAL BIOLOGY September 7, 2018

Dana C. Crawford, PhD Associate Professor Population and Quantitative Health Sciences Institute for Computational Biology



Single nucleotide variants (SNVs) and disease

3,541 publications69,969 SNP-trait associations

INSTITUTE FOR

COMPUTATIONAL BIOLOGY



Genome-wide arrays

~700,000 -millions variants

~\$40 - \$140/sample

Limited to known variants

Mostly common variants

European bias*



Single nucleotide variants (SNVs) and disease

Cardioprotection and APOC3 splice site, nonsense, and missense mutations MAF <1%



Sequencing

Assays known and unknown variants

Millions to billions of bases

Common, rare, and everything in between

~\$750 - \$1,500/sample



The Evolution of Sequencing

First generation sequencing (mid-1970s)

"Chemical sequencing"

- Add radioactive P³² to 5' phosphate of DNA
- Add chemical treatments to selectively remove DNA bases
- Run on gel to infer position









First generation sequencing (1977)

- "Chain Termination"
- Incorporate radio- or fluorescently labeled dideoxynucleotide triphosphates (ddNTPs)
- ddNTPs terminate polymerization
- Run on gel to infer position







ATGCAGCGTTACCATG...



First generation sequencing (1977-present)



COMPUTATIONAL BIOLOGY

Radioactive Gel-based One base per lane

> Fluorescent Gel-based One sample per lane



First generation sequencing (1977-present)



ABI 370A (1986)

Photo credit: Giac83 on wikimedia

INSTITUTE FOR COMPUTATIONAL BIOLOGY



ABI 310 (1995)

Photo credit: Jacopo Werther on wikimedia



ABI 3700 (1999)

Springer (2006) American Laboratory News

The Evolution of Sequencing

First generation sequencing (1977-present)

Throughput improved but limited by

- Separation methods
- Ability to parallel sequencing

Stratton, Campbell, Futreal (2009) *Nature* 458:719-724

INSTITUTE FOR

COMPUTATIONAL BIOLOGY



Next generation sequencing (2005-present)

Sequencing by synthesis (SBS):

Tracks bases being added to growing DNA strand

Pyrosequencing (Roche 454)

- Uses natural nucleotides (dATP, dTTP, etc)
- Detection based on pyrophosphate released during polymerase reaction





day per machine

Next generation sequencing (2005-present)

Throughput improved greatly improved, but

- Reads were short (80-120 bases) •
- Less accurate than Sanger (repeats, etc.)
- Need paired-end for *de novo* sequencing

Stratton, Campbell, Futreal (2009) Nature 458:719-724





Machine output (Mb)

Next generation sequencing (2005-present)

Trends in next-gen

- Longer reads
- Accuracy
- Price \$\$\$\$\$

Reuter, Spacek, Snyder (2015) Molecular Cell 58:586-597





Illumina HiSeq X Ten (2014)

\$1,000 genome

1.8 TB in <3 days

18K genomes/year

Genomes only



Illumina NovaSeq 6000 (2017)

\$100 genome?

80 GB-6 TB in 1-2 days

3-48 whole genomes/run

Genomes, exomes, RNA-Seq, ChIP-Seq, etc 2-Channel system (2 dyes)



https://www.ecseq.com/support/ngs/do_you_have_two_colors_or_four_colors_in_Illumina



PRECISION MEDICINE RESEARCH SEQUENCING

Enables population-scale whole genome sequencing for research

INSTITUTE FOR

COMPUTATIONAL BIOLOGY



Human Genomes Sequenced Annually

https://www.illumina.com/documents/products/illumina _sequencing_introduction.pdf

PRECISION MEDICINE RESEARCH



PRECISION MEDICINE RESEARCH SEQUENCING AlogUS



Aims to ascertain

- At least 1 million US residents
- Representative of the population

Willing to provide

- Biospecimens
- Health data (EHRs)
- Lifestyle/behavior data

Longitudinally

(Photo credit: Dr. Janina Jeff at the 2017 New Balance Bronx 10 Mile)

INSTITUTE FOR COMPUTATIONAL BIOLOGY



DATA FROM PRECISION MEDICINE RESEARCH SEQUENCING PERSISTENT BIAS Der the parts wen years, the proportion of participants in genome-wide asociation studies (GWAS) that are of Asian ancestry has increased.

A catalog of single nucleotide variants (SNVs) in diverse populations





DATA FROM PRECISION MEDICINE RESEARCH SEQUENCING

Variant servers with SNVs from diverse populations

Now includes eMERGEseq and PGRNseq with >24,000 samples and 1,138 genes



Sequence and Phenotype Integration Exchange (SPHINX) is a web-based tool for exploring data for hypothesis generation, especially around drug response implications of genetic variation across the eMERGE PGx cohort.



Bush et al (2016) Clin Pharmacol Ther 100:160-169

INSTITUTE FOR

COMPUTATIONAL BIOLOGY

DATA FROM PRECISION MEDICINE RESEARCH SEQUENCING

Variant servers with SNVs from diverse populations



NHLBI Exome Sequencing Project (ESP)

gnomAD browser beta | genome Aggregation Database

>6,000 sar

http://evs.

Search for a gene or variant or region

Example - Gene: PCSK9, Variant: 1-55516888-G-GA

>603000cample exome sequences \$150000xachblegeinstitetseoug/ces http://gnomad.broadinstitute.org/

Chr. Region:	1:100000-1100000
Single Chr. Location:	7:5567417
rsiD:	rs71531321



DATA FROM PRECISION MEDICINE

NHLBI Trans-Omics for Precision Medicine G Whole Genome Sequencing Program



Powered by TOPMed Freeze5 on GRCh38

(This dataset includes 463 million variants on 62784 individuals. Click here to switch to Freeze3a on GRCh37/hg19.)



https://bravo.sph.umich.edu/freeze5/hg38/



DATA FROM PRECISION MEDICINE

Variant: 11:116830637 C / T (APOC3: p.Arg19Ter, p.Arg37Ter)

Annotations Frequencies Sequence Depth Genotype Quality Raw Sequences Site Quality Metrics Summary

Summary		A	Annotations		
Filter Status Existing Variation Allele Frequency Allele Count Homozygous Alt Count CADD UCSC ClinVar	PASS rs76353203 0.0008362 105 / 125568 0 32.0 11-116830637-C-T C 11.116830637 is not in ClinVar C Open rs76353203 in ClinVar C	2	This variant falls on 5 transcripts belonging stop gained • APOC3 • ENST00000227667 (p.Arg19Ter) LOF: High-confidence • ENST00000433777 (p.Arg19Ter) LOF: High-confidence • ENST00000375345 (p.Arg37Ter)	to 1 gene: splice region • <i>APOC3</i> - ENST00000470144 (n.87C>T)	Exam va
			 ENST00000630701 (p.Arg37Ter) 		

LoF: High-confidence

💽 This list may not include additional transcripts in the same gene that the variant does not overlap.

INSTITUTE FOR COMPUTATIONAL BIOLOGY

ple TOPMed riant data

Population	Allele Frequency
1000G African	Not available
1000G American	Not available
1000G East Asian	Not available
1000G European	Not available
1000G South Asian	Not available
TOPMed Freeze5 on GRCh38	0.0008362

Frequency Table

SEQUENCING AND PRECISION MEDICINE RESEARCH EXAMPLES

General genotype-phenotype studies have been somewhat disappointing

Rare variant analysis will require

- Large (>10K) sample size
- "Careful" statistical analyses



LETTER

doi:10.1038/nature13917

Exome sequencing identifies rare *LDLR* and *APOA5* alleles conferring risk for myocardial infarction



SEQUENCING AND PRECISION MEDICINE RESEARCH EXAMPLES

General genotype-phenotype studies have been somewhat disappointing Natarajan et al (2018) Nat Comm 9:3391

"At these sample sizes and for these phenotypes, the incremental value of WGS for discovery is limited but WGS permits simultaneous assessment of monogenic and polygenic models to severe hypercholesterolemia."

INSTITUTE FOR

COMPUTATIONAL BIOLOGY



PRECISION MEDICINE RESEARCH



SEQUENCING

<u>??????</u>

(Photo credit: Dr. Janina Jeff at the 2017 New Balance Bronx 10 Mile)

INSTITUTE FOR COMPUTATIONAL BIOLOGY



SEQUENCING AND PRECISION MEDICINE RESEARCH CHALLENGES

Bioinformatics

- Pipelines and expertise
- Functional annotation
- Statistical methods

- Storage, secure access, and retrieval





Wilkins, Stokes, Wilson (1953) Nature 171:738-740



Franklin and Gosling (1953) Nature 171:740-741

QUESTIONS?



Watson and Crick (1953) *Nature* 171:737-738



SEQUENCING AND PRECISION MEDICINE RESEARCH EXAMPLES

Population studies and candidate genes



DATA FROM PRECISION MEDICINE RESEARCH SEQUENCING

A catalog of single nucleotide variants (SNVs)

RESEARCH ARTICLE

HUMAN GENETICS

Distribution and clinical impact of functional variants in 50,726 whole-exome sequences from the DiscovEHR study



DATA FROM PRECISION MEDICINE RESEARCH SEQUENCING

SPECIAL ARTICLE

INSTITUTE FOR

COMPUTATIONAL BIOLOGY

Genetic Misdiagnoses and the Potential for Health Disparities

Arjun K. Manrai, Ph.D., Birgit H. Funke, Ph.D., Heidi L. Rehm, Ph.D., Morten S. Olesen, Ph.D., Bradley A. Maron, M.D., Peter Szolovits, Ph.D., David M. Margulies, M.D., Joseph Loscalzo, M.D., Ph.D., and Isaac S. Kohane, M.D., Ph.D.

N Engl J Med 2016; 375:655-665 August 18, 2016 DOI: 10.1056/NEJMsa1507092

REPORT

Association of Trypanolytic ApoL1 Variants with Kidney Disease in African Americans

Giulio Genovese^{1,2,*}, David J. Friedman^{1,3,*}, Michael D. Ross⁴, Laurence Lecordier⁵, Pierrick Uzureau⁵, Barry I. Freedman⁶, Donald W. Bowden^{7,8}, Carl D. Langefeld^{8,9}, Taras K. Oleksyk¹⁰, Andrea L. Uscinski Knob⁴, Andrea J. Bernhardy¹, Pamela J. Hicks^{7,8}, George W. Nelson¹¹, Benoit Vanhollebeke⁵, Cheryl A. Winkler¹², Jeffrey B. Kopp¹¹, Etienne Pays^{5,†}, Martin R. Pollak^{1,13,†}

+ Author Affiliations

¹To whom correspondence should be addressed. E-mail: mpollak@bidmc.harvard.edu (M.R.P.); epays@ulb.ac.be (E.P.)

₽* These authors contributed equally to this work.

Science 13 Aug 2010: Vol. 329, Issue 5993, pp. 841-845 DOI: 10.1126/science.1193032