

# What data is being generated? Implications for returning results

Matt Lebo, PhD, FACMG

Chief Laboratory Director - Lab for Molecular Medicine
Director of Bioinformatics – Mass General Brigham Personalized Medicine
Associate Professor of Pathology – Brigham and Women's Hospital and Harvard Medical School
Associate Member – Broad Institute of MIT and Harvard

### Genomic assays

#### Genotyping array

- Typically target 500K 2 million **sites** across the genome
- Can identify other potential variant sites via imputation

#### Targeted panel

- Typically target 10s-100s of genes, though some are in low 1000s
- Can be enhanced to target specific clinically relevant content (e.g., CNVs or deep intronic variants)

#### Exome

Targets the ~2% of genome that is coding (i.e., all ~20,000-25,000 genes)

#### Genome

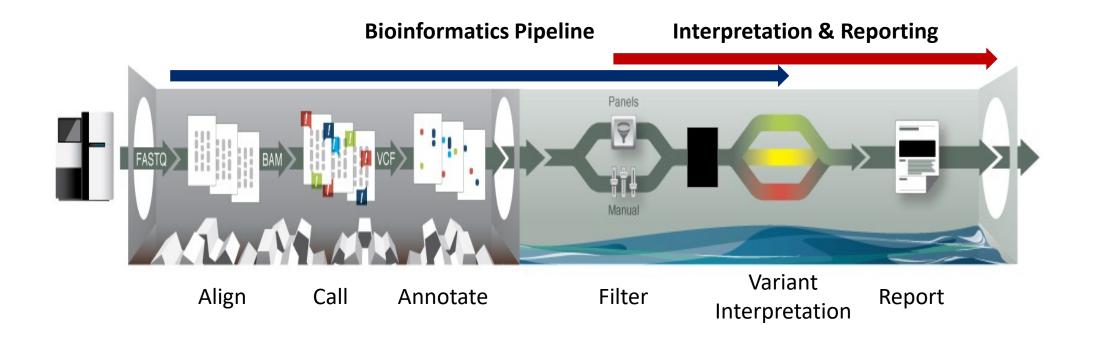
Unbiased sequencing of entire genome, including mitochondrial genome

#### Low-pass genome

Average coverage varies based upon costs and needs (0.25x to 5x)



Requires imputation to identify haplotypes and variants



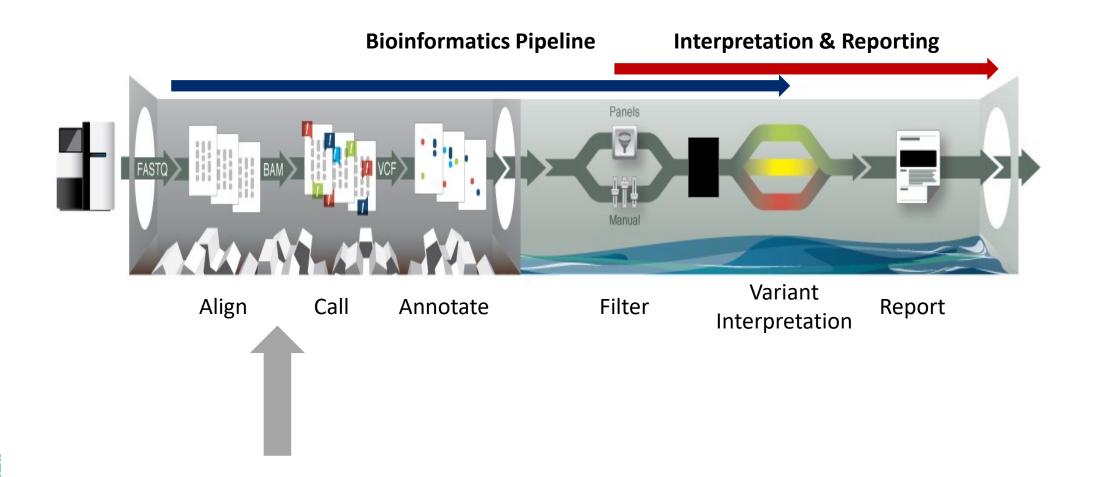


### Raw data and storage

	IDAT	CRAM	VCF	Imputed VCF
Genotyping array	60M	NA	20M	750M
Targeted Panel	NA	500M	5M	NA
Exome	NA	5G	100M	NA
Genome	NA	30G	1G	NA

- File sizes may differ, particularly for VCF
  - Joint/merged VCF files may reduce average file size
  - Annotated VCF files will increase file sizes
  - VCF files and annotations may be represented efficiently in a database

\*All sizes are estimates and vary on sequencing depth, panel size, and annotations from variant caller





### Alignment/Variant calling

#### **Clustering (for arrays)**

• Typically, not or can't be done for rare variants

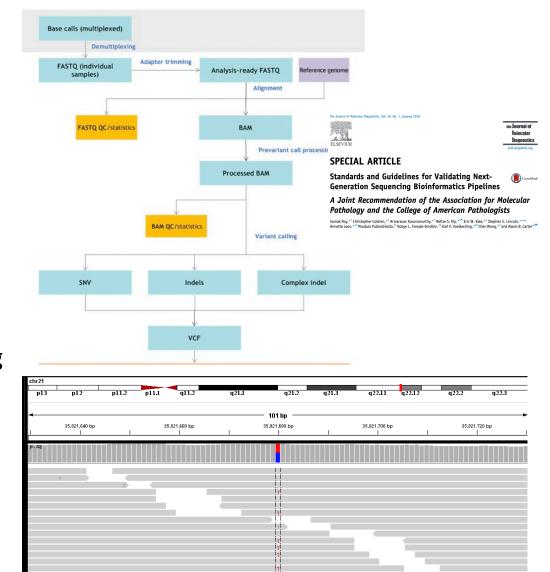
#### **Alignment**

- Choice of alignment tool affects accuracy
- Choice of reference genome may affect variant calling and downstream annotations

#### Variant calling

- Multiple variant callers may be needed to capture all relevant variant types
- Choice of caller affects accuracy





# Genomic assays – Relative accuracy

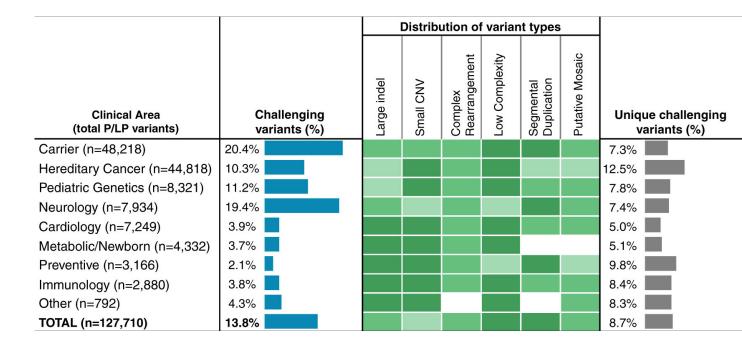
	Common variants	Rare Variants	Copy Number Variants	Structural Variants
Genotyping array	Good, but limited	Poor and limited	Good, but generally for large events	Poor
Targeted Panel	Good	Good	Low-resolution, unless specifically assayed	None
Exome	Good	Good	OK, low-resolution	Poor
Genome	Good	Good	Good, but low PPV for small events	Emerging
Low-pass genome	Good	Poor	Good for larger events	Unknown

<sup>\*</sup>All sequencing methodologies are using short reads



### Difficult regions (i.e., missing data)

- Even with generally high accuracy, many critical variants are still challenging to detect
- Important to know where datasets may be incomplete



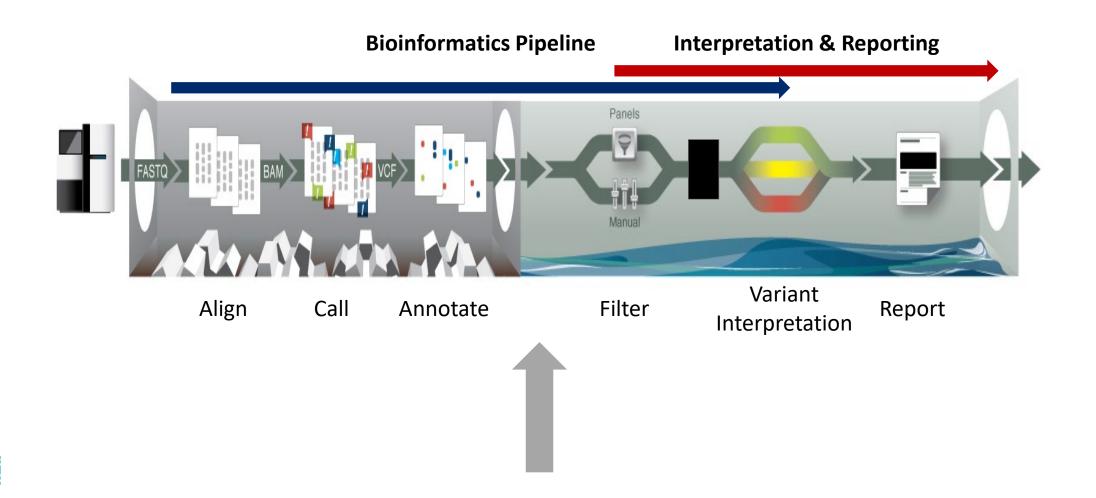
Genetics inMedicine

### www.nature.com/gim

#### ARTICLE

One in seven pathogenic variants can be challenging to detect by NGS: an analysis of 450,000 patients with implications for clinical sensitivity and genetic test implementation

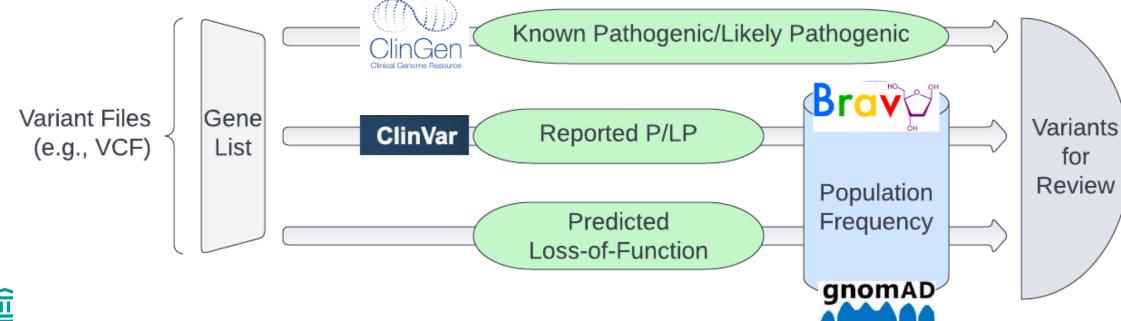






### Filtration Methods

- Need to annotate and filter variants
- Goal: identification of returnable Pathogenic and Likely Pathogenic variants
  - No prioritization based upon patient phenotype
- Balance sensitivity with PPV





### Standards in variant interpretation

Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology

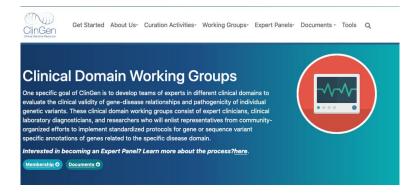
Sue Richards, PhD¹, Nazneen Aziz, PhD².¹6, Sherri Bale, PhD³, David Bick, MD⁴, Soma Das, PhD⁵,
 Julie Gastier-Foster, PhD⁵.³8, Wayne W. Grody, MD, PhD³.¹0.¹¹, Madhuri Hegde, PhD¹²,
 Elaine Lyon, PhD¹³, Elaine Spector, PhD¹⁴, Karl Voelkerding, MD¹³ and Heidi L. Rehm, PhD¹⁵;
 on behalf of the ACMG Laboratory Quality Assurance Committee

					SCORE:													Г
	Benign				I						tho	genic					_	
		Strong			Supporting	Ţ		Supporting			Moderate			Strong	_		Very Strong	Ξ
data	BAI	MAF too high								PM2	Absent (or rare) in pop db with coverage >20X		₩					
Population data	BSI	MAF too high				,	PS4_P	Proband Count - Supporting		PS4_M	Proband Count - Moderate		PS4	Case-control OR Proband Count				L
8	BS2	Observ in unaffected																
we data				BP1	Truncating disease causing; variant missense								PS1	Same AA change as establish pathogenic variant				
Computational and predictive data				BP3	In-frame indel in repeat region w/out known function					PM5	Diff pathogenic missense variant at codon		PM5_S	≥2 diff path missense variants at codon				
ntational ar				BP4	Computational evidence suggests no impact	F	PP3	Computational evidence suggests impact		PVS1_M	Null variant - Moderate		PVS1_S	Null variant - Strong		PVS1	Null variant & LOF known mechanism	
Сотр				BP7	Silent (or noncons splice) variant with no predicted splice impact					PM4	Protein length changing variant in non-repeat region							
Functional data						F	PP2	Missense in a gene with low rate of benign missense & path missense		PM1	Mutation hotspot or fxnl domain							
Funct	BS3	Estabished fxnl study shows no deleterious effect				F	PS3_P	Functional assay - Supporting		PS3_M	Functional assay - Moderate		PS3	Established fxnl study shows deleterious effect				
Seg Data	BS4	Lack of segregation in affected				1		Coseg with disease Dominant: 3 segs Recessive:		PP1_M	Coseg with disease Dominant: 5 segs Recessive:		PP1_S	Coseg with disease Dominant: 7 segs Recessive:				
De novo data										PM6	De novo (neither paternity or maternity confirmed)		PM6_S	≥2 independent occurences of PM6				
De nov													PS2	De novo (paternity and maternity confirmed)		PS2_VS	≥2 independent occurences of PS2	
Alleleic				BP2	Observed in trans with dominant variant OR observed in cis with path variant					РМ3	Detected in trans with path variant (recessive disorders)		PM3_S	2 occurences of PM3		PM3_VS	≥3 occurences of PM3	
Othe				BP6	ClinVar expert panel = benign	F	PP5	ClinVar expert panel = pathogenic										
Other data\				BP5	Found in case with an alternative cause	F	PP4	Patient phenotype or FH high specific for gene										Ĺ

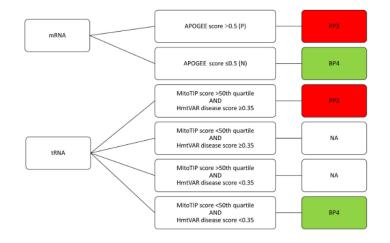
# Technical standards for the interpretation and reporting of constitutional copy-number variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics (ACMG) and the Clinical Genome Resource (ClinGen)

Erin Rooney Riggs, MS, CGC<sup>1</sup>, Erica F. Andersen, PhD<sup>2,3</sup>, Athena M. Cherry, PhD<sup>4</sup>, Sibel Kantarci, PhD<sup>5</sup>, Hutton Kearney, PhD<sup>6</sup>, Ankita Patel, PhD<sup>7</sup>, Gordana Raca, MD, PhD<sup>6</sup>, Deborah I. Ritter, PhD<sup>9</sup>, Sarah T. South, PhD<sup>10</sup>, Erik C. Thorland, PhD<sup>6</sup>, Daniel Pineda-Alvarez, MD<sup>11</sup>, Swaroop Aradhya, PhD<sup>4,11</sup> and Christa Lese Martin, PhD<sup>1</sup>

Section 1: Initial assessment of genomic cont	ent			
Evidence type	Evidence	Suggested points/case	Max	
			score	
Copy-number loss content	<ol> <li>Contains protein-coding or other known functionally important elements.</li> </ol>	0 (Continue evaluation)	0	
	<ol> <li>Does NOT contain protein-coding or any known functionally important elements.</li> </ol>	-0.60	-0.60	
	haploinsufficiency (HI) or established benign genes/genomic regions (Skip to section	3 if your copy-number loss DOES NOT	overlap the	
types of genes/regions)				
Overlap with ESTABLISHED HI genes or genomic	<ol> <li>Complete overlap of an established HI gene/genomic region.</li> </ol>	1.00	1.00	
regions and consideration of reason for referral				
	2B. Partial overlap of an established HI genomic region	0 (Continue evaluation)	0	
	<ul> <li>The observed CNV does NOT contain the known causative gene or critical region for this</li> </ul>			
	established HI genomic region OR			
	Unclear if known causative gene or critical region is affected OR			
	No specific causative gene or critical region has been established for this HI			
	genomic region			
	2C. Partial overlap with the 5' end of an established HI gene (3' end of the gene not	See categories below		
	involved)			
	2C-1and coding sequence is involved	0.90 (range: 0.45 to 1.00)	1.00	
	2C-2and only the 5' UTR is involved	0 (range: 0 to 0.45)	0.45	
	2D. Partial overlap with the 3' end of an established HI gene (5' end of the gene not involved)	See categories below		
	2D-1and only the 3' untranslated region is involved.	0 (Continue evaluation)	0	
	2D-2and only the last exon is involved. Other established pathogenic variants have	0.90 (range: 0.45 to 0.90)	0.90	
	been reported in this exon.			
	2D-3and only the last exon is involved. No other established pathogenic variants have	0.30 (range: 0 to 0.45)	0.45	
	been reported in this exon.			
	2D-4and it includes other exons in addition to the last exon. Nonsense-mediated	0.90 (range: 0.45 to 1.00)	1.00	
	decay is expected to occur.			
	2E. Both breakpoints are within the same gene (intragenic CNV; gene-level sequence	See ClinGen SVI working group	See	
	variant).	PVS1 specifications	categor	
		<ul> <li>PVS1 = 0.90</li> </ul>	at left	
		(Range: 0.45 to 0.90)		
		<ul> <li>PVS1_Strong = 0.45</li> </ul>		
		(Range: 0.30 to 0.90)		
		<ul> <li>PVS1_Moderate or PM4 (in-frame</li> </ul>		
		indels) = 0.30		
		(Range: 0.15 to 0.45)		
		<ul> <li>PVS1_Supporting = 0.15</li> </ul>		
		(Range: 0 to 0.30)		



### Specifications of the ACMG/AMP standards and guidelines for mitochondrial DNA variant interpretation

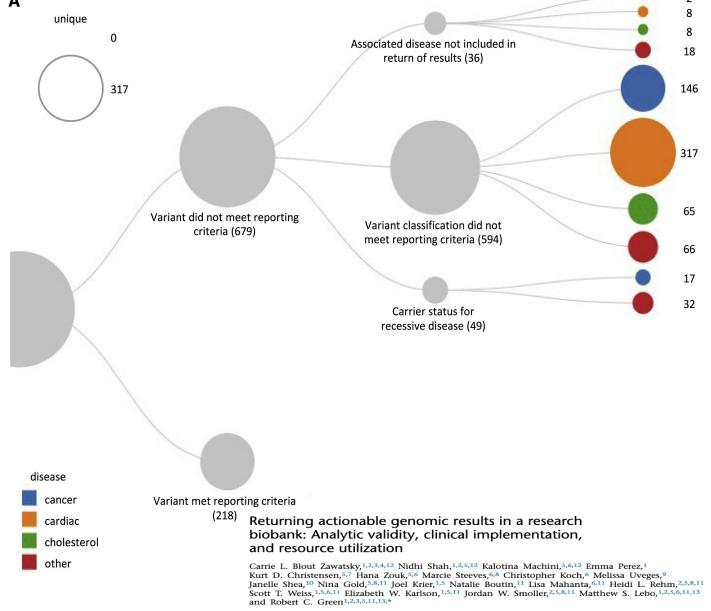




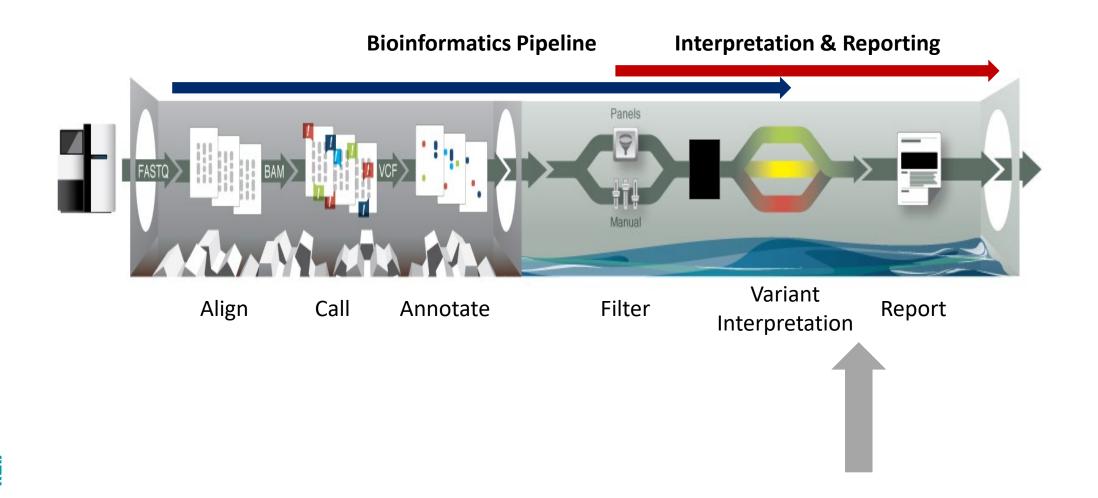
### **Filtration Results**

# Results still require manual review for pathogenicity

- Not all previously reported variants meet criteria for P/LP
- Not all variants annotated as loss-offunction are actually LOF
- Not all diseases associated to a gene are returnable
- Variants may only be in carrier state for recessive disease









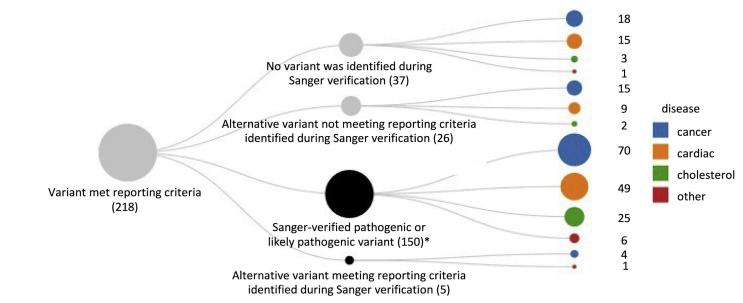
## Returnable results: genotyping arrays vs. sequencing

#### **Genotyping arrays**

- Predefined list of variants interrogated
- Poor performance for rare variants
  - Inaccurate sites are often recurrent
  - Can have TP and FP for same site

#### Sequencing

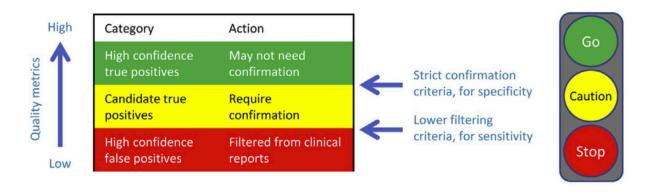
- Interrogates whole coding sequence
  - Can get "novel" variants
- General high accuracy of variant calling
  - However, confirmation may still be necessary for to identify false positive calls



Returning actionable genomic results in a research biobank: Analytic validity, clinical implementation, and resource utilization

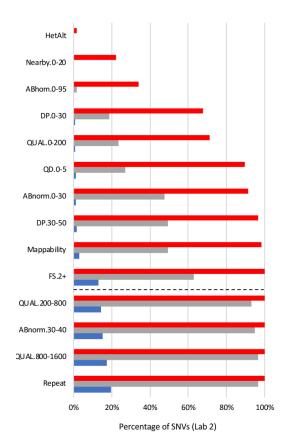
Carrie L. Blout Zawatsky,<sup>1,2,3,4,12</sup> Nidhi Shah,<sup>1,2,5,12</sup> Kalotina Machini,<sup>5,6,12</sup> Emma Perez,<sup>1</sup> Kurt D. Christensen,<sup>5,7</sup> Hana Zouk,<sup>5,6</sup> Marcie Steeves,<sup>6,8</sup> Christopher Koch,<sup>6</sup> Melissa Uveges,<sup>9</sup> Janelle Shea,<sup>10</sup> Nina Gold,<sup>5,8,11</sup> Joel Krier,<sup>1,5</sup> Natalie Boutin,<sup>11</sup> Lisa Mahanta,<sup>6,11</sup> Heidi L. Rehm,<sup>2,5,8,13</sup> Scott T. Weiss,<sup>1,5,6,11</sup> Elizabeth W. Karlson,<sup>1,5,11</sup> Jordan W. Smoller,<sup>2,5,8,11</sup> Matthew S. Lebo,<sup>1,2,5,6,11,13</sup> and Robert C. Green<sup>1,2,3,5,11,13,\*</sup>

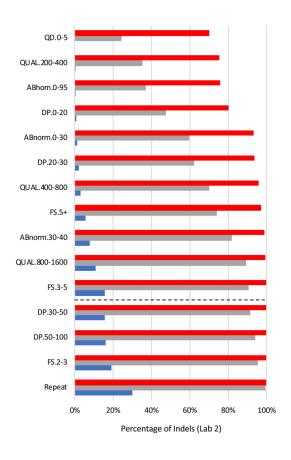
## Confirmation of sequencing variants



A Rigorous Interlaboratory Examination of the Need to Confirm Next-Generation
Sequencing—Detected Variants with an Orthogonal
Method in Clinical Genetic Testing

Stephen E. Lincoln,\* Rebecca Truty,\* Chiao-Feng Lin,†† Justin M. Zook,§ Joshua Paul,\* Vincent H. Ramey,\* Marc Salit,§¶
Heidi L. Rehm,†‡||\*\*†† Robert L. Nussbaum,\*†‡ and Matthew S. Lebo††\*\*††

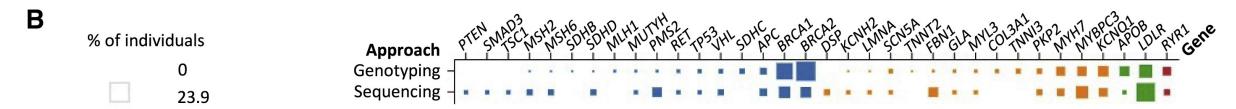




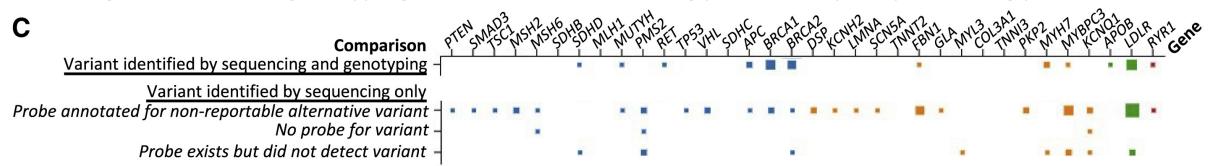


## Returnable results: genotyping arrays vs. sequencing

#### Sequencing more accurately reflects true spectrum of variation



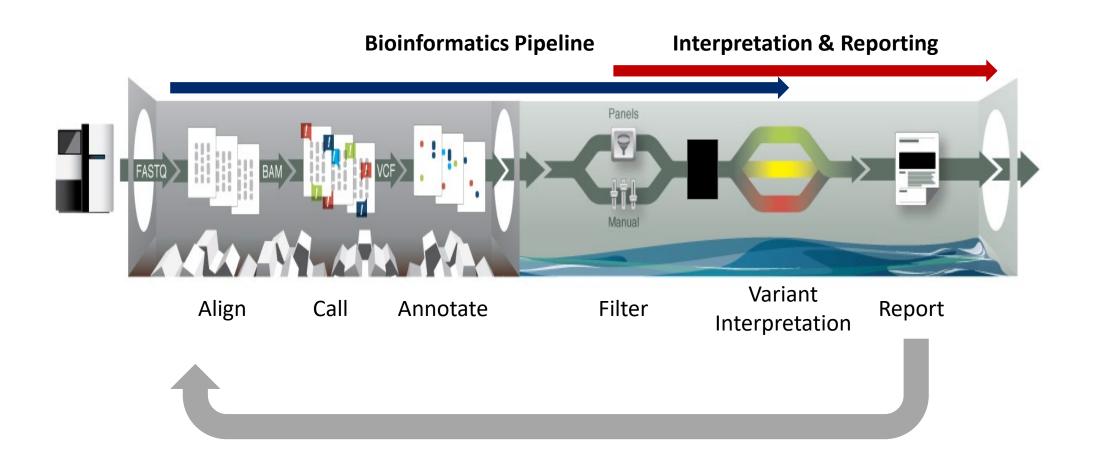
#### Missing variation from genotyping data often due to missing probes and poor performing probes



Returning actionable genomic results in a research biobank: Analytic validity, clinical implementation, and resource utilization

Carrie L. Blout Zawatsky,<sup>1,2,3,4,12</sup> Nidhi Shah,<sup>1,2,5,12</sup> Kalotina Machini,<sup>5,6,12</sup> Emma Perez,<sup>1</sup> Kurt D. Christensen,<sup>5,7</sup> Hana Zouk,<sup>5,6</sup> Marcie Steeves,<sup>6,8</sup> Christopher Koch,<sup>6</sup> Melissa Uveges,<sup>9</sup> Janelle Shea,<sup>10</sup> Nina Gold,<sup>5,8,11</sup> Joel Krier,<sup>1,5</sup> Natalie Boutin,<sup>11</sup> Lisa Mahanta,<sup>6,11</sup> Heidi L. Rehm,<sup>2,5,8,11</sup> Scott T. Weiss,<sup>1,5,6,11</sup> Elizabeth W. Karlson,<sup>1,5,11</sup> Jordan W. Smoller,<sup>2,5,8,11</sup> Matthew S. Lebo,<sup>1,2,5,6,11,13</sup> and Robert C. Green<sup>1,2,3,5,11,13,\*</sup>







### Revisiting raw data

#### Reassessment of variants can lead to additional returns

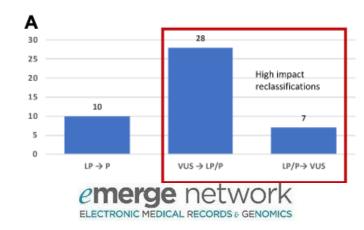
- New evidence available
- New guidelines for classification

#### Improved algorithms can more accurately call variants

- Especially true for CNVs/SVs
- Also for small variants, even with current high accuracy

#### Other improvements affecting bioinformatic pipelines

- Annotations (e.g., updated transcripts, LOF prediction)
- Improvements in reference genomes



	Variant type	Sensitivity	PPV
Pipeline 1	SNV	0.9900	0.9941
	Indel	0.9897	0.9937
Pipeline 2	SNV	0.9967	0.9993
	Indel	0.9948	0.9967



Reanalysis of eMERGE phase III sequence variants in 10,500 participants and infrastructure to support the automated return or knowledge updates

Hana Zouk • Wanfeng Yu • Andrea Oza • ... Scott T. Weiss • Matthew S. Lebo • Heidi L. Rehm 2 20 • School of the subsection of the support of the s

Published: November 30, 2021 • DOI: https://doi.org/10.1016/j.gim.2021.10.010 • 🌘

### Improvements in the human genome reference

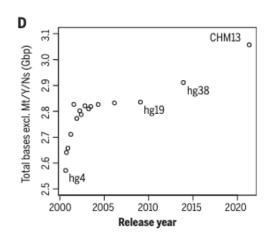
- Telomere-To-Telomere
  - Complete haploid genome (CHM13)
  - Additional 200Mb of genomic content

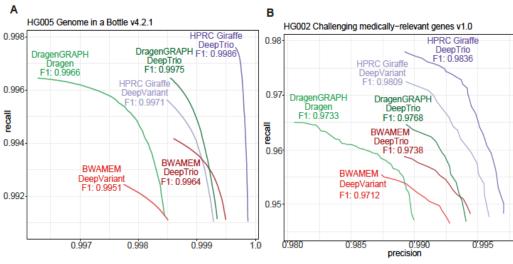
- Human Pangenome Reference
  - Captures genetic diversity of human species
  - Improved variant calling, especially in difficult regions
  - Improved SV/CNV calling

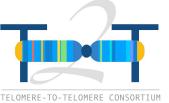
#### A Draft Human Pangenome Reference

📵 Wen-Wei Liao, Mobin Asri, Iana Ebler, Daniel Doerr, Marina Haukness, 📵 Glenn Hickey, 📵 Shuangiia Lu, Julian K. Lucas, De Jean Monlong, Haley J. Abel, Silvia Buonaiuto, De Xian H. Chang, Haoyu Cheng, Justin Chu, Vincenza Colonna, D Jordan M. Eizenga, Xiaowen Feng, Christian Fischer, Robert S. Fulton, Shilpa Garg, Cristian Groza, Andrea Guarracino, William T Harvey, Simon Heumos, Kerstin Howe, Miten Jain, Tsung-Yu Lu, 📵 Charles Markello, 📵 Fergal J. Martin, Matthew W. Mitchell, 📵 Katherine M. Munson, Moses Njagi Mwaniki, O Adam M. Novak, D Hugh E. Olsen, D Trevor Pesout, D David Porubsky, D Piotr Prins, o Jonas A. Sibbesen, Chad Tomlinson, Flavia Villani, Mitchell R. Vollger, Human Pangenome Reference Consortium, @ Guillaume Bourque, @ Mark JP Chaisson, @ Paul Flicek, Adam M. Phillippy, Justin M. Zook, D. Evan E. Eichler, D. David Haussler, Erich D. Jarvis, D. Karen H. Miga, Ting Wang, 10 Erik Garrison, Tobias Marschall, 10 Ira Hall, 10 Heng Li, 10 Benedict Paten doi: https://doi.org/10.1101/2022.07.09.499321



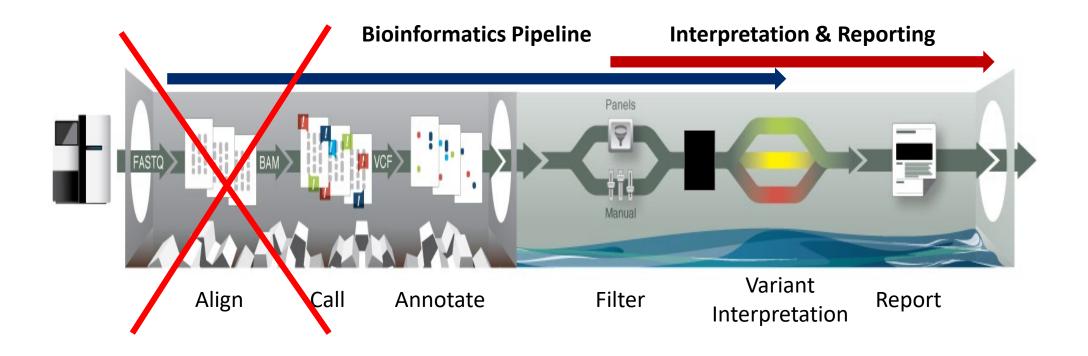






#### The complete sequence of a human genome

SERGEY NURK , SERGEY KOREN , ARANG RHIE , MIKKO RAUTIAINEN , ANDREY V. BZIKADZE , ALLA MIKHEENKO, MITCHELL R. VOLLGER NICOLAS ALTEMOSE D, LEV URALSKY D, [...], AND ADAM M. PHILLIPPY D +90 authors Authors Info & Affiliations SCIENCE - 31 Mar 2022 - Vol 376, Issue 6588 - pp. 44-53 - DOI: 10.1126/science.abj6987





### Bioinformatic considerations for NOT reporting

#### Can you generate data without initiating return?

- Most cohorts have not masked/removed potentially returnable variants
  - Typically release unannotated VCF and/or aggregate data
  - Difficult to identify actionable variants without further information

#### Are there methods to mask rare variants?

- For genotyping arrays, you can remove them from the manifest
  - May appear in imputation if "common" enough
- Harder for sequencing-based methods
  - Due to identification of "novel" variants
  - Possible to do it by population frequency (e.g., only keep variants >5% frequency)
    - Will remove more variants than may be wanted



# Thank you!

