

Big Data: Genomic Reference Databases to Empower Mendelian Diagnosis

Anne O'Donnell-Luria, MD, PhD

Associate Director for Rare Disease Genomics

Broad Institute of MIT and Harvard

Clinical Geneticist, Boston Children's Hospital

Twitter: AnneOtation





- NIH-funded center launched in early 2016 to discover new disease-gene relationships underlying <u>Mendelian</u> disease
- We work with collaborators with existing cohorts of patient samples consented for genetic studies, prescreened for some known causes of disease



- CMG covers cost of exome sequencing; supports analysis
- Diagnoses & gene discoveries are pursued and published by collaborator
- Commitment to data sharing



Mendelian Inheritance: Use this to choose one of the built-in inheritance search methods. Variants will be returned that segregate with the selected inheritance model. Methods are described here.

- O Homozygous Recessive
- X-Linked Recessive
- Compound Heterozygous
- O Dominant
- De Novo Dominant





seqr analysis software



Bamshad et al., Nature Reviews Genetics (2011) 12, 745-755.





Map against reference genome Determine variants, Filter, compare patients



Clinical exome sequencing in a new tool in our diagnostic tool box

- Sequence ~20,000 human genes
- 10,000 30,000 protein coding variants



What's in an ex

- Every genome contains many rare, potentially functional variants
 - ~500 rare missense variants (1/3 of wh predicted damaging by *in silico* predict
 - ~100 LoF variants: ~20 homozygous, ~20 rare
 - o ~100 rare variants in known disease genes
 - ~50 reported disease-causing mutations (!)
 - o 1-2 *de novo* coding mutations
 - 0

How can we identify the pathogenic variant(s) in the sea of benign vari





Harnessing the power of allele frequency



Making sense of one exome requires **tens of thousands** of exomes (or genomes) to reveal rare variants



Five-fold reduction in number of very rare variants with large reference databases

- # variants remaining in an exome after applying a 0.1% filter across all populations
- Both size and ancestral diversity increase filtering power

Lek et al., Nature, 2016





http://www.internationalgenome.org/1000-genomes-browsers/



http://evs.gs.washington.edu/EVS



http://exac.broadinstitute.org/





http://gnomad.broadinstitute.org/

The genome aggregation database (gnomAD)

- Data provided by 107 PIs for >138,000 individuals including 123,136 exomes & 15,496 whole genomes
- Illumina data, processed through same pipeline, called jointly
- Sites VCF of entire dataset available for download -> Can annotate your dataset with allele frequencies
- Individual level data not shared & phenotype data not available
- Cases and controls from common disease studies. No Mendelian disease studies knowingly included.
- New population (e.g. >5K Ashkenazi Jewish samples)
- Report the population with the highest allele frequency for each variant (popmax AF)
- 55% Male; Mean age 54 years

<u>http://gnomad.broadinstitute.org</u> <u>http://gnomad-beta.broadinstitute.org</u>



- African (12,942)
- Latino (18,237)
- Ashkenazi Jewish (5,081)
- East Asian (9,472)
- Finnish European (13,046)
- European (63,416)
- South Asian (15,450)

Ancestry and sex are inferred f principal component analysis (rather than self-reported

Sample QC Removes

Low quality samples Sex chromosome abnormalit First and second degree relati



gnomAD browser beta | genome Aggregation Database

CFTR

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

About gnomAD

The Genome Aggregation Database (gnomAD) is a resource developed by an international coalition of investigators, with the goal of aggregating and harmonizing both exome and genome sequencing data from a wide variety of large-scale sequencing projects, and making summary data available for the wider scientific community.

The data set provided on this website spans 123,136 exome sequences and 15,496 whole-genome sequences from unrelated individuals sequenced as part of various disease-specific and population genetic studies. The gnomAD Principal Investigators and groups that have contributed data to the current release are listed here.

Recent News

January 25, 2018

Beta release of redesigned gene and region pages.

October 3, 2017

gnomAD r2.0.2 released. Sample composition is identical to the previous release (r2.0.1), however we have made a change to the variant filtering







Matthew Solomonson Nick Watts

Konrad Karczewski Ben Weisburd http://gnomad.broadinstitute.org http://gnomad-beta.broadinstitute.org

CFTR cystic fibrosis transmembrane conductance regulator (ATP-binding cassette sub-family C, member 7)



CFTR cystic fibrosis transmembrane conductance regulator (ATP-binding cassette sub-family C, member 7)



gnomAD variant page

CFTR Phe508del chr7:117199644 ATCT / A

Variant: 7:117199644 ATCT / A

	Exomes	Genomes	Total			
Filter	Pass	Pasa				
Allele Count	1712	235	1947			
Allele Number	246048	30934	276982			
Allele Frequency	0.006958	0.007597	0.007029			
dbSNP	rs1801178					
UCSC	7-117199644-ATCT-A 🖸					
ClinVar	Click to sear	ch for variant i	n Clinvar 🖓			

Genotype Quality Metrics

Site Quality Metrics

tote: This variant is multiallelic! The other alt alleles are:

Report this variant

7-117199644-A-G

Annotations

This variant falls on 5 transcripts in 2 genes:

inframe deletion



Note: This list may not include additional transcripts in the same gene that the variant does not overlap

Population -	Allele Count	Allele Number	Number of *	Allele Frequency
European (Non- Finnish)	1531	126568	1	0.01210
Ashkenazi Jewish*	59	10150	0	0.005813
Other	36	6454	0	0.005578
Latino	133	34412	0	0.003865
African	64	24024	0	0.002664
European (Finnish)	61	25738	0	0.002370
South Asian	63	30776	0	0.002047
East Asian	0	18860	0	0.000
Total	1947	276982	1	0.007029

* For detailed analysis of Ashkenazi Jewish frequency see the IBD Exomes Browser.

Read Data

This interactive IGV.js visualization shows reads that went into calling this variant.

Note: These are reassembled reads produced by GATK HaplotypeCaller --bamOutput so they accurately represent what HaplotypeCaller was seeing when it called this variant.



Raw read data supporting a variant is available

http://gnomad-beta.broadinstitute.org/variant/7-117199644-ATCT-A



gnomAD variant page

CFTR Phe508del chr7:117199644 ATCT / A

Expect to see **9 homozygotes** in 63,000 Europeans

- Carrier frequency as predicted
- Severe pediatric-onset disease cases depleted (but not entirely removed)

Population Frequencies

Population	Allele Count	Allele ♦ Number	Number ofHomozygotes	 Allele Frequency
European (Non- Finnish)	1531	126568	1	0.01210
Ashkenazi Jewish*	59	10150	0	0.005813
Other	36	6454	0	0.005578
Latino	133	34412	0	0.003865
African	64	24024	0	0.002664
European (Finnish)	61	25738	0	0.002370
South Asian	63	30776	0	0.002047
East Asian	0	18860	0	0.000
Total	1947	276982	1	0.007029

Do you think the homozygote is a real variant? - Review the read data



Considerations for gnomAD IGV visualization of variants

- Low confidence loss of function (LC LOF)
- Poorly aligned regions (ex: low copy repeat)
- Multinucleotide variants (MNVs)
- Homopolymer runs
- Complex indels
- Somatic mosaicism

Low confidence loss of function variants

- LOFTEE flags variants that are unlikely to cause loss of function, for example:
 - Dubious transcript annotation
 - Protein truncating variant near end of the gene

10:81319089 T / A (rs549800979)	E	G p.Lys51Ter	stop gained		11
10:81319140 CGGGGATACCA / C (rs767274650)	E	p.Gly31AlafsTer70	frameshift		1
10:81319171 CT / C (rs779122123)	E	G p.Lys23ArgfsTer81	frameshift		5
10:81319920 C / T		G c218-1G>A†	splice acceptor	LC LoF	3
10:81320119 C / A (rs150972292)		G c54+1G>T	splice donor	LC LoF	1
10:81320119 C / T (rs150972292)		G c54+1G>A	splice donor	LC LoF	65

Poorly aligned regions

• Multiple variants in region

• Different allele balances

• Raises concern about variants called in this region



Poorly aligned regions

• Multiple variants in region

• Different allele balances

9:35906557 C / A (rs747119413)

(rs747119413)

9:35906558 C / A

9:35906558 C / CA

9:35906557 CCACCCCGCCA / C

9:35906559 A / AC (rs781316793)

9:35906559 A / C (rs781316793)

 Raises concern about variants called in this region

E

E

E

E

G



Homopolymer runs

- Homopolymer G
 - Indels in these regions enriched for PCR artifacts
 - But also region enriched for true variants



Multinucleotide variants

- Two variants within 1 codon in vcf considered separately but should be interpreted together
- Multinucleotide variants (MNV)
 - Variant 1: T>C, Ser>Pro (missense)
 - Variant 2: C>A, Ser>* (nonsense)
 - MNP: TC>CA, Ser>Gln (missense)
- These are flagged in ExAC, working on them for gnomAD
- Can see similar situation with complex indels (deletion and insertion that maintain the frame

ACTGTTTCAAGAATTCCACAAA Sequence Coverage Paired-end reads

Somatic mosaicism

- See skewed allele balance
- Many of these are filtered but not all



When a variant is absent from gnomAD, it's important to determine if that region is covered

gnomAD browser beta		About	Downloads	Terms	Contact	Jobs	FAQ
	Interested in working on the development of this resource? Apply here.						

gnomAD browser beta | genome Aggregation Database

6:1611497

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

Unable to find variant in gnomAD

Possible reasons:

1)This is not the position in the canonical transcript displayed on the browser
2)Position is not covered in gnomAD
3)Variant is not in gnomAD

Look up chromosome coordinate at http://mutalyzer.nl

Coverage summary

Missense + LoF

LoF

Looking for: chr6:1611497 C > A Pro273Thr

Look for the closest variant

Pro273Thr is not present but Pro273Pro is present

65K chromosomes or 32.5K people genotyped at this position

						Indels			
Export table to CSV									
Variant	- Source #	Consequence	Annotation	≑ Flags		Allele Number	* Number of Homozygotes	Allele Freque	ency ¢
6:1611478 C / A	G	p.Ser266Arg	missense		1	29296	0	3.413e-5	
6:1611485 A / G (rs747812123)	E	p.Ser269Gly	missense		1	74862	0	1.336e-5	
6:1611486 G / A	EG	p.Ser269Asn	missense		3	106734	0	2.811e-5	
6:1611487 C / A	E	p.Ser269Arg	missense		0	77082	0	0	
6:1611487 C / G	E	p.Ser269Arg	missense		4	106184	0	3.767e-5	
6:1611488 G / C	EG	p.Gly270Arg	missense		5	106352	0	4.701e-5	
6:1611492 G / T	E	p.Ser271lle	missense		0	73722	0	0	
6:1611495 G / GC	E	p.Ser276GInfsTer30	frameshift	LoF flag	2	67846	0	2.948e-5	
6:1611495 GC / G	E	p.Pro274ArgfsTer41	frameshift	LoF flag	3	67846	0	4.422e-5	
6:1611496 C / T	E	p.Ser272Ser	synonymous		0	67854	0	0	
6:1611499 C / T		p.Pro273Pro	synonymous		1	65354	0	1.53e-5	
6:1611500 C / T (rs769224835)	E	p.Pro274Ser	missense		1	65280	0	1.532e-5	
6:1611513 C / T (rs777329712)	E	p.Pro278Leu	missense		1	50484	0	1.981e-5	



Include:

Exomes

Conomo

SNPs

Jadala I

Filtered (non-PASS) variants

Evaluating rare variant pathogenicity

© American College of Medical Genetics and Genomics

beh

vogenic

ACMG STANDARDS AND GUIDELINES

Genetics inMedicine 2015

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, PhD^{2,16}, Sherri Bale, PhD³, David Bick, MD⁴, Soma Das, PhD⁵, Julie Gastier-Foster, PhD^{6,7,8}, Wayne W. Grody, MD, PhD^{9,10,11}, 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

		Benign		Pathogenic			
_		Strong	Supporting	Supporting	Moderate	Strong	Very strong
	Population data	MAF is too high for disorder BA1/BS1 OR observation in controls inconsistent with disease penetrance BS2			Absent in population databases PM2	Prevalence in affecteds statistically increased over controls PS4	
	Computational and predictive data		Multiple lines of computational evidence suggest no impact on gene /gene product BP4 Missense in gene where only truncating cause disease BP1 Silent variant with non predicted splice impact BP7 In-frame indels in repeat w/out known function BP3	Multiple lines of computational evidence support a deleterious effect on the gene /gene product PP3	Novel missense change at an amino acid residue where a different pathogenic missense change has been seen before PM5 Protein length changing variant PM4	Same amino acid change as an established pathogenic variant PS1	Predicted null variant in a gene where LOF is a known mechanism of disease PVS1
	Functional data	Well-established functional studies show no deleterious effect BS3		Missense in gene with low rate of benign missense variants and path. missenses common PP2	Mutational hot spot or well-studied functional domain without benign variation PM1	Well-established functional studies show a deleterious effect PS3	
	Segregation data	Nonsegregation with disease BS4		Cosegregation with disease in multiple affected family members PP1	Increased segregation data	>	
Richards et al.,	De novo data				De novo (without paternity & maternity confirmed) PM6	De novo (paternity and maternity confirmed) PS2	
Genet Med, 2015	Allelic data		Observed in <i>trans</i> with a dominant variant BP2 Observed in <i>cis</i> with a pathogenic variant BP2		For recessive disorders, detected in trans with a pathogenic variant PM3		

Identification of constrained genes


Identification of constrained genes in ExAC





pLI identifies known haploinsufficient genes for pediatric-onset conditions

JAG1

Constraint from ExAC	Expected no. variants	Observed no. variants	Constraint Metric	
Synonymous	235.6	200	z = 1.44	
Missense	478.1	297	Z = 4.05	
LoF	45.2	1	pLI = 1.00	

Alagille syndrome (dominant congenital disorder affecting liver, heart and eyes)

Probability of loss-of-function (LOF) intolerance: pLI scores

- Haploinsufficiency: in a diploid organism, where having only one functional copy of a gene is insufficient to sustain a wild type phenotype and leads to a "abnormal" phenotype.
- pLI >0.9 is considered evidence of haploinsufficiency
- 3,230 genes have pLI score >0.9
- 70% have not been assigned a phenotype in OMIM
- We predict that loss of function variation in these genes will result in disease or embryonic lethality

pLI does **not** identify genes haploinsufficient genes for <u>adult-onset</u> conditions

Majority of disease impact is post-fertility

BRCA	1 Constraint from ExAC	Expected no. variants	Observed no. variants	Constraint Metric
	Synonymous	204.2	210	z = -0.25
	Missense	508.9	567	z = -1.26
	LoF	46.1	36	pLI = 0.00

Breast and ovarian cancer

pLI does not identify genes for recessive conditions

Constraint from ExAC	Expected no. variants	Observed no. variants	Constraint Metric	
Synonymous	171.6	168	z = 0.17	
Missense	418.8	671	z = -6.03	
LoF	53.2	54	pLI = 0.00	

Cystic fibrosis (recessive disorder affecting lungs and pancreas)

Missense constraint

JAG1

Constraint from ExAC	Expected no. variants	Observed no. variants	Constraint Metric
Synonymous	235.6	200	z = 1.44
Missense	478.1	297	Z = 4.05
LoF	45.2	1	pLI = 1.00

Missense constrained genes have a Z-score > 3

~1800 missense constrained genes

Regional missense constraint https://www.biorxiv.org/content/early/2017/06/12/148353

ClinVar has a growing catalog of variant interpretations but VUSes remain a major challenge



Zach Zappala

Constraint on the browser

Haploinsufficiency results in Kleefstra syndrome



http://gnomad-beta.broadinstitute.org/gene/KMT2C

Constraint on the browser



http://gnomad-beta.broadinstitute.org/gene/KMT2C

Constraint on the browser



http://gnomad-beta.broadinstitute.org/gene/KMT2C

Gene expression on the browser



http://gnomad-beta.broadinstitute.org/gene/KMT2C

Gene expression on the browser





Gene constraint @

ExAC | gnomAD

	Exp. no. Ot	os. no.	Constraint
/	Median across all tissues (5.22)		metric
	Specific tissue		Z = -0.40
	Adipose subcutaneous (4.945)		
	Adipose visceral omentum (4.7)		Z = 1.53
	Adrenal gland (4.61)		pLI = 1.00
	Artery aorta (5.19)		
	Artery coronary (5.22)		
	Artery tibial (6.78)		
	Bladder (6.28)		
	Brain amygdala (4.445)		
	Brain anteriorcingulatecortex ba24 (3.	99)	
	Brain caudate basalganglia (4.3)		
	Brain cerebellarhemisphere (9.91)		Isoform expression @
	Brain cerebellum (10.08)		Median across all tissue
	Brain cortex (4.02)		
	Brain frontalcortex ba9 (4.3)		•
	Brain hippocampus (4.275)		•
	Brain hypothalamus (5.095)	(4.00)	•
	Brain nucleusaccumbens basalganglia	(4.82)	•
	Brain putamen basalganglia (4.205)		•
	Brain spinalcord cervical c1 (5.46)		
	_		

1.42

1.23

0.34

http://gnomad-beta.broadinstitute.org/gene/KMT2C

Also on http://exac.broadinstitute.org

0.93

GTEx Portal



Example gene expression across tissues: KMT2C



https://www.gtexportal.org/home/gene/KMT2C

Example gene expression across tissues: KMT2C



https://www.gtexportal.org/home/gene/KMT2C

Can a male with OTC be an organ donor?

- Adult male with OTC deficiency (urea cycle defect) presented with brain herniation in the setting of illness
- Declared brain dead
- Had requested organ donation are organs from someone with OTC deficiency safe for transplantation?
- Literature review
 - Has been done successfully except for liver (would result in OTC deficiency in recipient)
 - Case reports of deaths in cases of undiagnosed OTC carrier females were liver donors
- Ask the experts
- Look at expression in different tissues

OTC Gene Expression



ASXL1 additional sex combs like 1 (Drosophila)

Gene constraint @

ExAC | gnomAD



Bohring-Opitz syndrome (BOS): Severe dominant disorder caused by protein truncating variants (PTVs) in *ASXL1*



- Well-established severe autosomal dominant pediatric-onset disorder
- Profound intellectual disability & characteristic facial features
- We would not expect to see any individuals with this disorder in ExAC or gnomAD

Collaboration with <u>University of Utah/ARUP</u> Colleen Carlston Hunter Underhill Tatiana Tvrdik Rong Mao

Clinical exome sequencing result: *De novo* dominant *ASXL1* p.R404* nonsense pathogenic variant

Population	Allele Count	\$ Allele Number ≑	Number of Homozygotes	\$ Allele Frequency
European (Finnish)	1	6614	0	0.0001512
East Asian	1	8654	0	0.0001156
African	1	10404	0	9.612e-05
European (Non-Finnish)	4	66710	0	5.996e-05
Latino	0	11578	0	0
Other	0	908	0	0
South Asian	0	16510	0	0
Total	7	121378	0	5.767e-05



- Are there patients with Borhing-Opitz syndrome in ExAC? **No**
- Does this variant cause Borhing-Opitz syndrome? Yes

chr20:31021211 C>T

There are numerous PTVs in ASXL1 in ExAC



PTVs found in ExAC, excluding individuals from the TCGA cohort

Carlston*, O'Donnell-Luria*, et al., Hum Mut, 2017

Read support for ASXL1 p.R404* shows skewed allele balance





Most ExAC ASXL1 PTVs show skewed allele balance



Carlston*, O'Donnell-Luria*, et al., Hum Mut, 2017

ExAC PTVs in *ASXL1* show skewed allele balance compared to other rare variants in *ASXL1*



Allele balance percentage

Clonal Hematopoiesis of Indeterminate Potential (CHIP)

- A well described phenomenon of aging
- Somatic mutations in certain genes provide a growth advantage to hematopoietic stem cells
- ASXL1 PTVs are known driver mutations in hematopoietic cancer



Increase in the risk of all-cause mortality, highest for hematologic cancer but also for solid tumors, coronary heart disease & ischemic stroke.



If PTVs in ExAC are due to clonal hematopoiesis of indeterminate potential (CHIP), then the *ASXL1* PTVs should be seen at higher frequency with increasing age



- This is consistent with ExAC ASXL1 PTVs arising by somatic mosaicism and clonal expansion, so are <u>not</u> germline.
- The **germline** p.R404* variant is pathogenic in the patient for BOS.

Carlston*, O'Donnell-Luria*, et al., Hum Mut, 2017

We can learn interesting biology from reference population databases starting from a single variant and a clinical question

Frequency filtering









James Ware Imperial College London

Nicky Whiffin Imperial College London

Eric Minikel HMS/Broad

Daniel MacArthur Broad/MGH/HMS

Central tenet

• The frequency of a pathogenic variant in a reference sample, *that is not selected for the condition*, should not exceed the prevalence of the condition.

Possible Exceptions

- Founder mutations and Bottlenecked populations
- Balancing selection
- Penetrance needs to be considered

Disease specific allele frequency (AF) thresholds for autosomal dominant disease

Genetic architecture



Whiffin*, Minikel* *et al*. Genetics in Medicine (2017)

Hypertrophic cardiomyopathy (HCM) specific AF threshold

Genetic architecture

0.5 x 1/500

disease prevalence X heterogeneity

penetrance

50%

3%

Most common pathogenic allele

MYBPC3:c.1504C>T causes 2.2% (1.6-3.0%) of European HCM cases

maximum credible population allele frequency

6x10⁻⁵

Whiffin*, Minikel* *et al*. Genetics in Medicine (2017)

Online calculator: cardiodb.org/alleleFrequencyApp

Frequency Filter HOME

calculate AF calculate AC

explore architecture inverse AF

penetrance about



James Ware

Allele frequencies: not exactly what they appear to be



Reference population databases sample the general population so we need to apply statistical estimates of uncertainty.

We have the ability to estimate the upper limit on the CI.

Lek et al., Nature, 2016

Precomputed across 5 ExAC populations: Filtering AF



- Allele count (AC) at the upper bound of the one-tailed 95% CI
- Specified as the maximum credible AF given the sample
- Computed for 5 main populations (AFR, AMR, EAS,
- Highest filtering AF reported

Rarity necessary but not sufficient for pathogenicity

Example: Looking up a ClinVar VUS

ClinVar Entry

NM_000256.3(MYBPC3):c.961G>A (p.Val321Met)							
Variation ID: 🕜 Review status: 🕜	161310 🗙 🚖 🚖 criteria provided, multiple submitters, no conflicts						
Interpretation 🕢	Interpretation 🕢						
Clinical significance: Last evaluated: Number of submission(s):	<u>Uncertain significance</u> Nov 8, 2016 9						

Variant: 11:47367887 C / T

Every variant page in ExAC has a Filtering AF (Coming soon for gnomAD)

MYBPC3

Filter Status dbSNP Allele Frequency	PASS rs200119454	Filtering all If the variant is too comm	ele frequency (AF): a threshold for filtering AF is greater than the man non to be causative and may be t	r filtering variants that are too common to plausibly aximum credible population AF for the disease of int filtered. Click here to see the filtering AF calculator a	cause disease. erest, the variant pp and citation.
Filtering AF	0.0007 (European	n (Non-Finnish))		Site Quality Metrics	
UCSC ClinVar	11-47367887-C-T Click to search for v	7 /ariant in Clinvar (Z		

missense	PW4 transcripts in	r genes:		Population	Allele Count	Allele Number	Number of Homozygotes ≑	Allele Frequency
• MYBPC3	Transcripts -			European (Non- Finnish)	36	45024	0	0.0007996
Note: This list may	ENST00000256993 (p.Val321Met)		A loes not overlap	African	1	6832	0	0.0001464
internet internet	ENOTO000000			East Asian	0	5714	0	0
	Polyphen: b	100000399249 (p.Val321Met) Volyphen: benign; SIFT: deleterious		European (Finnish)	0	4320	0	0
	ENST0000544	791 (p.Val321Met)		Latino	0	6552	0	0
	Polyphen: possibly_damaging; SIFT: dele		1	Other	0	584	0	0
	ENST0000545	ENST00000545968 * (p.Val321Met)		South Asian	0	10974	0	0
	Polyphen: benign; SIFI: deleterious			Total	37	80000	0	0.0004625

Population Frequencies
Filtering AF for HCM

Filtering AF 0.0007 (European (Non-Finnish))



Filtering AF for HCM

Filtering AF 0.000593 (European (Non-Finnish))



Classification by ACMG criteria

MYBPC3 c.961G>A, p.Val321Met

Using frequency filter approach, we can say: BS1 too common in controls

BUT still need to consider other evidence (if there is any)

Other criteria met: BP5 alternate cause found in several cases

No segregation data available No functional data available

Likely benign

Richards et al., Genetics in Medicine, 2015

	Ben	ign	Pathogenic						
	Strong	Supporting	Supporting	Moderate	Strong	Very strong			
Population data	MAF is too high for disorder BA1/BS1 OR observation in controls inconsistent with disease penetrance BS2			Absent in population databases PM2	Prevalence in affecteds statistically increased over controls PS4				
Computational and predictive data		Multiple lines of computational evidence suggest no impact on gene /gene product BP4 Missense in gene where only truncating cause disease BP1 Silent variant with non predicted splice impact BP7 In-frame indels in repeat wfout known function BP3	Multiple lines of computational evidence support a deleterious effect on the gene /gene product PP3	Novel missense change at an amino acid residue where a different pathogenic missense change has been seen before PM5 Protein length changing variant PM4	Same amino acid change as an established pathogenic variant PS1	Predicted null variant in a gene where LOF is a known mechanism of disease PVS1			
Functional data	Well-established functional studies show no deleterious effect BS3		Missense in gene with low rate of benign missense variants and path. missenses common PP2	Mutational hot spot or well-studied functional domain without benign variation PM1	Well-established functional studies show a deleterious effect PS3				
Segregation data	Nonsegregation with disease BS4		Cosegregation with disease in multiple affected family members PP1	Increased segregation data	→				
De novo data				De novo (without paternity & maternity confirmed) PM6	De novo (paternity and maternity confirmed) PS2				
Allelic data		Observed in trans with a dominant variant BP2 Observed in cis with a pathogenic variant BP2		For recessive disorders, detected in trans with a pathogenic variant PM3					
Other database		Reputable source w/out shared data = benign BP6	Reputable source = pathogenic PP5						
Other data		Found in case with an alternate cause BP5	Patient's phenotype or FH highly specific for gene PP4						



Confidence:

121412

○ 0.9 • 0.95 ○ 0.99 ○ 0.999

Reference population size (alleles)

Maximum credible population AF:

1e-05

- Dominant, complete penetrance
- Prevalence 1:10,000
- Only 1 gene causes phenotype
- Most common pathogenic variant accounts for 20% of cases

Maximum tolerated reference AC:



Maximum credible population AF:

2e-05

- Dominant, **50% penetrance**
- Prevalence 1:10,000
- Only 1 gene causes phenotype
- Most common pathogenic variant accounts for 20% of cases



Maximum tolerated reference AC:



Maximum credible population AF:

0.002

- Recessive, fully penetrant
- Only 1 gene causes phenotype
- Most common pathogenic variant accounts for 20% of cases



Maximum tolerated reference AC:

269



ClinGen disease expert panel working groups drafting guidelines for disease specific allele frequency thresholds

https://www.clinicalgenome.org/

Conclusions

- Reference population databases are critically important to evaluate variant rarity, which is necessary but not sufficient for pathogenicity for rare disease
- Constrained genes show less variation among humans than expected and are enriched for genes that result in disease when mutated
- Frequency filtering is a more stringent, statistically-based approach to set allele frequency cut offs for variant filtering and interpretation
- The power of reference population datasets will increase as they grow in size and diversity
 - gnomAD v3 with ~60,000 genomes anticipated by early 2019

Frequently asked questions

- Phenotypes: Very limited phenotype information and regulatory restrictions on sharing – need for phenotypegenotype databases (biobanks)
- Subsets: Non-cancer, non-neuro coming with next gnomAD release in Fall 2018
- Constraint on gnomAD: Coming with next gnomAD release
 - Genes that do not have constraint mainly annotation issue (Gencode) or too many variants (synonymous and missense) often related to mapping issues (pseudogenes)

Click here to contribute data	Geno ₂ MP						
to Geno ₂ MP	PCGF2	Search					
HPO Browser Need to find the HPO term for a clinical finding?	 Gene (MYH3); chromosome position of variant (17:10534960); dbSNP rsID (rs34393601) HPO term (oral cleft); HPO number (0000202) 						
	Rare variants (<1% AF) from ~8,	,000 samples					
Gene PC	CGF2,CISD3						
Description	on polycomb group ring finger 2;CDGSH iron sulfur domain 3						
Number of variants	ts 95						
GeneCard							
MalaCards	ds PCGF2 CISD3						
OMIN	M PCGF2 CISD3						

Gene summary

Filter by annotation category All											Export table to TSV			
Chr:Pos \$	Alleles	rsID	HPO Profiles \$\\$	# het 😄	# hom ≑	Gene 🌲	mRNA 💠	Annotations \$\\$	cDNA Change	Protein Change	ESP AC 💠	ExAC AC 💠	1K Genome AC 💠	CADD Phred-scaled \$\\$
17:36894794	T>C	NA	2	3	0	PCGF2	NM_007144.2	synonymous-near	c.480A>G	p.(K160=)	0	0	0	6.016
17:36890664	G>A	NA	1	0	1	PCGF2,CISD3	NM_001136498.1	3-prime-UTR	c.*956G>A	NA	0	0	0	3.901
17:36891006	G>T	NA	1	0	1	PCGF2,CISD3	NM_001136498.1	3-prime-UTR	c.*1298G>T	NA	0	0	0	1.246
17:36891329	CAT>C	rs71764170	5	3	2	CISD3,PCGF2	NM_001136498.1	3-prime-UTR	c.*1622_*1623del2	NA	0	0	0	-1
17:36891402	G>A	rs1061140	2	2	1	PCGF2,CISD3	NM_001136498.1	3-prime-UTR	c.*1694G>A	NA	0	0	0.005952	14.45
17:36891443	G>C	rs376471012	1	1	0	PCGF2,CISD3	NM_001136498.1	3-prime-UTR	c.*1735G>C	NA	0.000085	0	0	10.39
17:36891457	A>T	rs141585937	1	1	0	PCGF2,CISD3	NM_001136498.1	3-prime-UTR	c.*1749A>T	NA	0.010982	0.014342	0.008242	7.646
17:36891467	T>A	rs377749543	1	1	0	PCGF2,CISD3	NM_001136498.1	3-prime-UTR	c.*1759T>A	NA	0.00008	0.000177	0	10.57
17:36891486	G>C	NA	2	2	0	PCGF2,CISD3	NM_001136498.1	3-prime-UTR	c.*1778G>C	NA	0	0	0	16.94
17:36891494	G>A	rs112908318	3	4	0	PCGF2,CISD3	NM_001136498.1	3-prime-UTR	c.*1786G>A	NA	0	0.000302	0.001374	10.92
17:36891495	G>A	NA	2	2	0	PCGF2,CISD3	NM_001136498.1	3-prime-UTR	c.*1787G>A	NA	0	0.000149	0	23.2
17:36891500	G>A	NA	1	1	0	PCGF2,CISD3	NM_001136498.1	3-prime-UTR	c.*1792G>A	NA	0	0.000146	0	9.086
17:36891506	G>A	NA	7	9	0	PCGF2,CISD3	NM_001136498.1	3-prime-UTR	c.*1798G>A	NA	0	0.00042	0	10.56
17:36891566	C>T	rs149121439	1	2	0	PCGF2,CISD3	NM_001136498.1	3-prime-UTR	c.*1858C>T	NA	0.000078	0	0	12.82

http://geno2mp.gs.washington.edu/

ClinVar PCGF2 CISD3

Variant seen in patients with nervous system abnormalities but also in unaffected relatives

Geno₂MP Phenotypes represented Hover over the bars to see number of individuals. Affected Relative blood & blood formi... genitourinary system nervous system connective lissue head and neck immune system integument tebolismihomeo... musculature respiratory system ardiovascular syst... neoplasm skeletal system growth ear

Variant count by HPO Profile

All

Filter by Sample Status

Total number of HPO profiles: 5

Export table to TSV

C/G	Sample Status	÷	Contact	HPO term: broad	HPO ID: broad	HPO term: medium	HPO ID: medium	HPO term: narrow	HPO ID: narrow
1	affected		<u>×</u>	Abnormality of the nervous system	HP:0000707	Behavioral abnormality	HP:0000708	Autism	HP:0000717
				Abnormality of the nervous system	HP:0000707	Seizures	HP:0001250		
1	affected		×	Abnormality of the nervous system	HP:0000707	Abnormality of the cerebellum	HP:0001317		
1	affected		<u>×</u>	Abnormality of the nervous system	HP:0000707	Neurodevelopmental abnormality	HP:0012759		
1	relative		×	Abnormality of the nervous system	HP:0000707	Abnormality of the cerebellum	HP:0001317	Cerebellar atrophy	HP:0001272
1	relative		\mathbf{x}	Abnormality of the nervous system	HP:0000707	Abnormality of the cerebellum	HP:0001317	Cerebellar hypoplasia	HP:0001321

There is power in big data when deployed in publicly available, intuitive user interfaces

Thank you to all the groups that contribute data to gnomAD and other public resources

Acknowledgements

Email: odonnell@broadinstitute.org

UCL: Nicky Whiffin, James Ware





Konrad Karczewski



Grace Francioli Tiao



Beryl Cummings

Daniel

MacArthur



Matthew Solomonson

Nick Watts

Alföldi

Monkol Lek

gnomAD Pls

gnomad.broadinstitute.org/about

Broad Data Sciences Platform

Eric Banks **Charlotte Tolonen** Christopher Llanwarne Dave Shiga Fengmei Zhao Jeff Gentry Jose Soto Kathleen Tibbetts Khalid Shakir **Kristian Cibulskis Miguel Covarrubias** Ryan Poplin Ruchi Munshi Sam Novod Thibault Jeandet Valentin Ruano-Rubio Yossi Farjoun Intel GenomicsDB **Google Cloud**

team (QC)

Cotton Seed Tim Poterba Jonathan Bloom Jacqueline Goldstein Dan King Ben Neale

hai



Kaitlin Samocha



Qingbo Wang

Eric

Minikel





Laura Gauthier



Kristen

Laricchia



Mike Wilson



Mark

Daly

NIGMS R01 GM104371 NIDDK U54 DK105566 **Broad Institute**







Collaboration with

University of Utah/ARUP

Colleen Carlston

Hunter Underhill

Tatiana Tvrdik

Rong Mao

Jessica

Publicly available reference population databases



Publicly available reference population databases



https://discovehrshare.com/

Publicly available reference population databases



https://search.hli.io

Mutational model accurately predicts synonymous variation

 We used our mutational model to predict the expected number of variants in the ~61K individuals in ExAC



Samocha et al., Nat Genetics, 2014; Lek et al., Nature, 2016