

Use of the Genome Aggregation Database (gnomAD)

Anne O'Donnell-Luria, MD, PhD

Broad Institute of MIT and Harvard

Boston Children's Hospital

H3Africa Rare Disease Workshop February 2021

Every genome contains many rare, potentially functional variants

- ~2 million high quality variants in a variant call file (vcf)
- ~20,000 variants within coding (exome) regions
- ~500 rare missense variants (1/3 of which are predicted

In Mendelian disease analysis, how can we identify the pathogenic genetic variant(s) in the sea of benign variation?

- ~50 reported disease-causing mutations (!)
- 1-2 *de novo* coding mutations
- Unknown number of sequencing errors

The power of 7.7 billion people

Given known mutation rates, it is almost certain that
**every possible single base change compatible with
life exists in a living human**

Variant aggregation at Broad

Exome Aggregation

Consortium (ExAC) – v1

60,076 exomes

Genome Build 37

released October 2014

Genome Aggregation

Database (gnomAD) – v2

125,748 exomes + 15,708 genomes

Genome Build 37

released October 2016

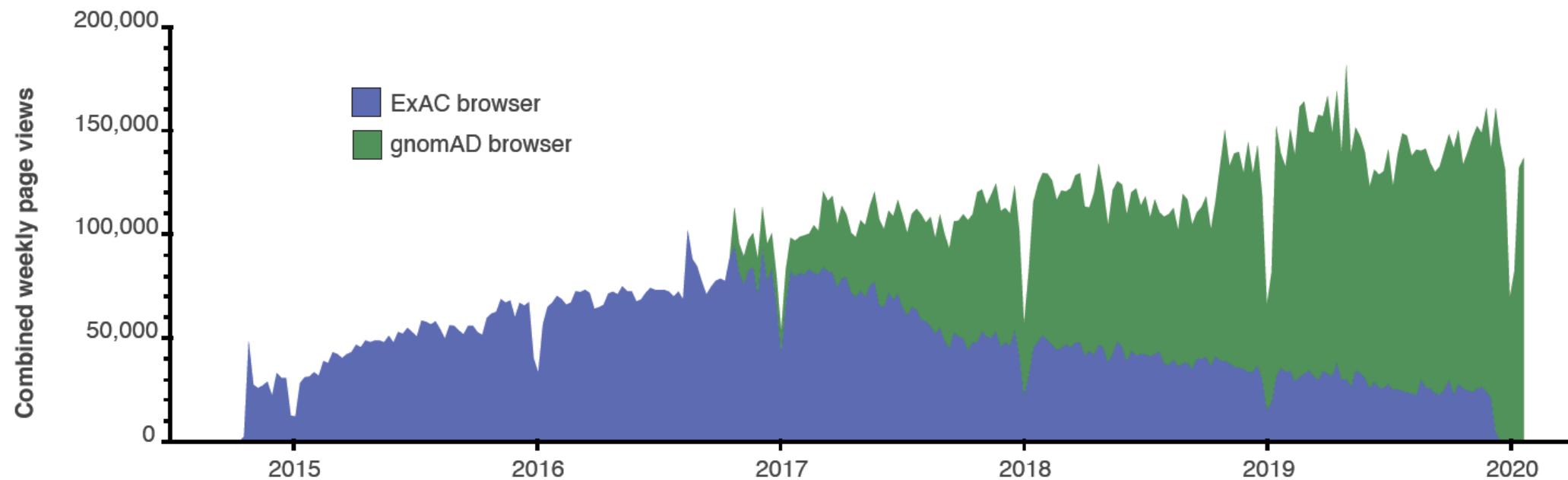
Genome Aggregation

Database (gnomAD) – v3

71,702 genomes

Genome Build 38

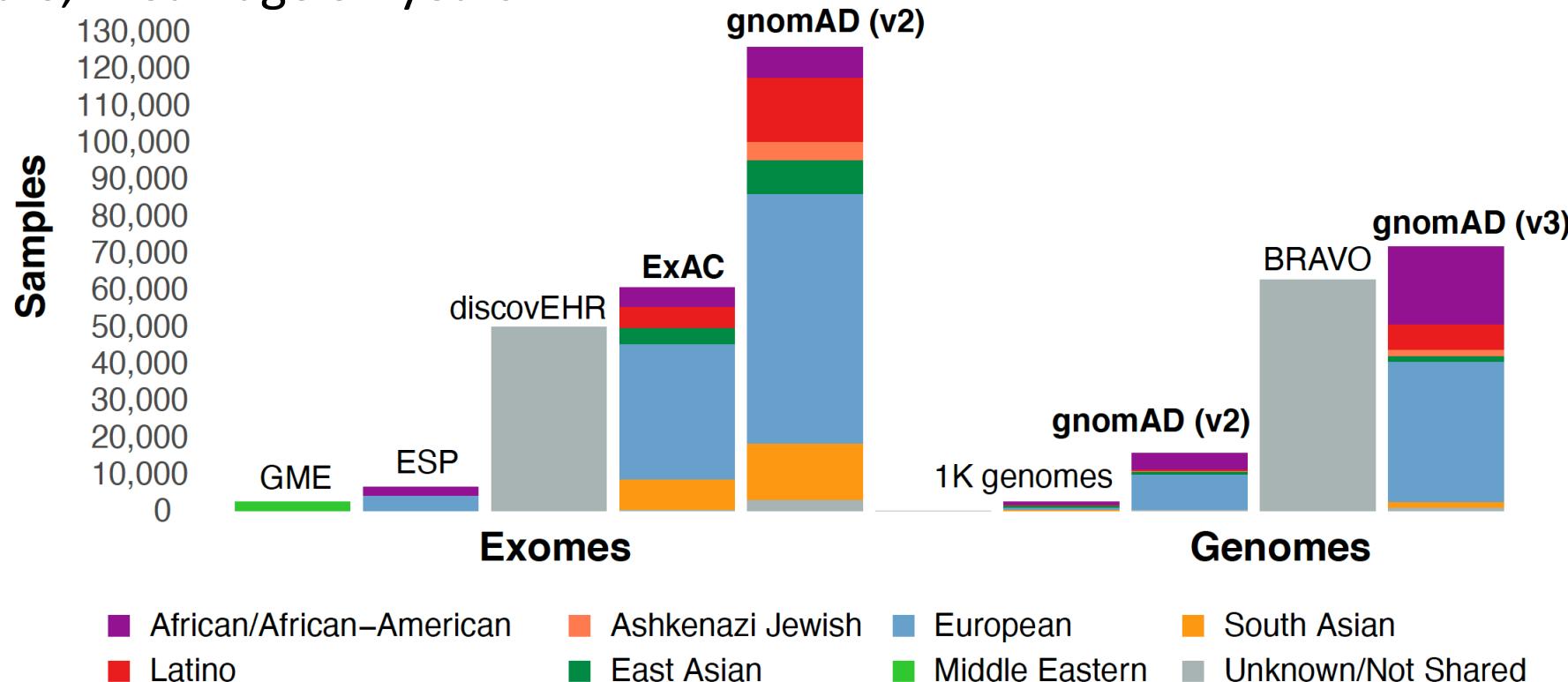
released October 2019



- 20+ M pageviews from 184 countries
- Aided in the diagnosis of over 200,000 patients with rare disease

Who's in the gnomAD database?

- Case-control studies of complex adult-onset diseases (e.g. type 2 diabetes, heart attack, migraine, bipolar)
- Depleted as much as possible of people **known** to have severe pediatric disease, and their first-degree relatives
- 55% Male, Mean age 54 years



Thank you to the PIs who share data in gnomAD

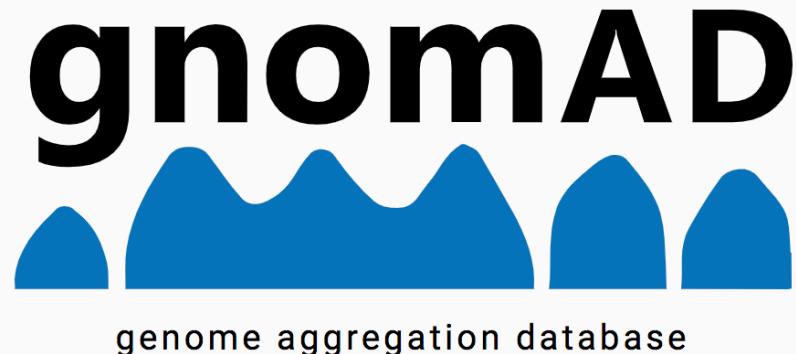
Data processed through same pipeline and called jointly

Maria Abreu	Wendy Chung	Ira Hall	Jaume Marrugat	Nazneen Rahman	Hilkka Soininen
Carlos A Aguilar Salinas	Bruce Cohen	Craig Hanis	Kari Mattila	Alex Reiner	Harry Sokol
Tariq Ahmad	Adolfo Correa	Matthew Harms	Steve McCarroll	Anne Remes	Tim Spector
Christine M. Albert	Dana Dabelea	Mikko Hiltunen	Mark McCarthy	Stephen Rich	Nathan Stitzel
Nicholette Allred	Mark Daly	Matti Holi	Jacob McCauley	John D. Rioux	Patrick Sullivan
David Altshuler	John Danesh	Christina Hultman	Dermot McGovern	Samuli Ripatti	Jaana Suvisaari
Diego Ardissino	Dawood Darbar	Chaim Jalas	Ruth McPherson	Dan Roden	Kent Taylor
Gil Atzman	Joshua Denny	Mikko Kallela	James Meigs	Jerome I. Rotter	Yik Ying Teo
John Barnard	Ravindranath Duggirala	Jaakko Kaprio	Olle Melander	Danish Saleheen	Tuomi Tiinamaija
Laurent Beaugerie	Josée Dupuis	Sekar Kathiresan	Deborah Meyers	Veikko Salomaa	Ming Tsuang
Gary Beecham	Patrick T. Ellinor	Eimear Kenny	Lili Milani	Nilesh Samani	Dan Turner
Emelia J. Benjamin	Roberto Elosua	Bong-Jo Kim	Braxton Mitchell	Jeremiah Scharf	Teresa Tusie Luna
Michael Boehnke	Jeanette Erdmann	Young Jin Kim	Aliya Naheed	Heribert Schunkert	Erkki Vartiainen
Lori Bonnycastle	Tõnu Esko	George Kirov	Saman Nazarian	Svati Shah	James Ware
Erwin Bottinger	Martt Färkkilä	Jaspal Kooner	Ben Neale	Moore Shoemaker	Hugh Watkins
Donald Bowden	Diane Fatkin	Seppo Koskinen	Peter Nilsson	Tai Shyong	Rinse Weersma
Matthew Bown	Jose Florez	Harlan M. Krumholz	Michael O'Donovan	Edwin K. Silverman	Maija Wessman
Steven Brant	Andre Franke	Subra Kugathasan	Dost Ongur	Pamela Sklar	James Wilson
Hannia Campos	Gad Getz	Soo Heon Kwak	Lorena Orozco	Gustav Smith	Ramnik Xavier
John Chambers	David Glahn	Markku Laakso	Michael Owen		
Juliana Chan	Ben Glaser	Terho Lehtimäki	Colin Palmer		
Daniel Chasman	Stephen Glatt	Ruth Loos	Aarno Palotie		
Rex Chisholm	David Goldstein	Steven A. Lubitz	Kyong Soo Park		
Judy Cho	Cicerio Gonzalez	Ronald Ma	Carlos Pato		
Rajiv Chowdhury	Leif Groop	Daniel MacArthur	Ann Pulver		
Mina Chung	Christopher Haiman	Gregory M. Marcus	Dan Rader		



Broad Genomics Platform
Broad Data Sciences Platform

This is a new version of the gnomAD browser. The old version is available at <http://gnomad-old.broadinstitute.org>



Search by gene, region, or variant

Examples - Gene: [PCSK9](#), Variant: [1-55516888-G-GA](#)

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.



<http://gnomad.broadinstitute.org>

MPG Primers on using gnomAD

https://www.youtube.com/watch?v=J29_Pf8Ti1s

<https://www.youtube.com/watch?v=Bh8AKkl-DhY>

NSD1 nuclear receptor binding SET domain protein 1

Dataset gnomAD v2.1.1 gnomAD SVs v2.1 ?

Genome build GRCh37 / hg19

Ensembl gene ID ENSG00000165671.14

Ensembl canonical transcript ENST00000439151.2

Other transcripts ENST00000347982.4, ENST00000503056.1, and 13 more

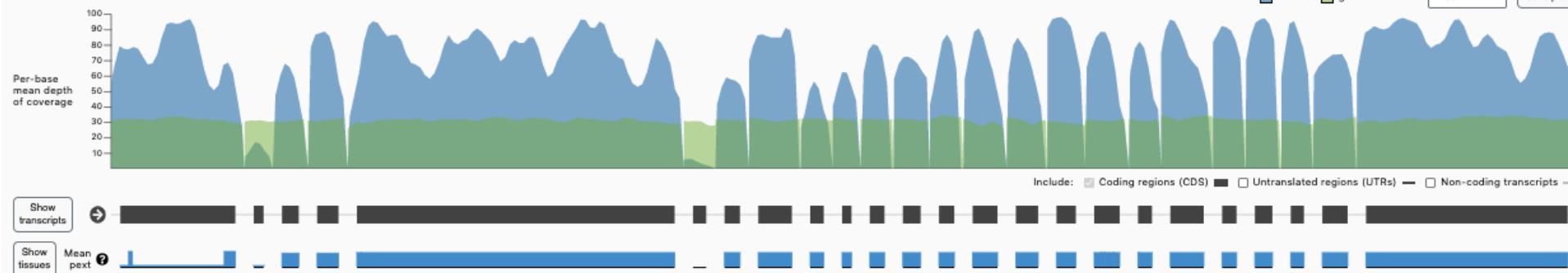
Region 5:176560026-176727216

References Ensembl, UCSC Browser, and more

Constraint ?

Category	Exp. SNVs	Obs. SNVs	Constraint metrics
Synonymous	525.3	523	Z = 0.08 o/e = 1 (0.93 - 1.07) 0 → 1
Missense	1427.7	1065	Z = 3.41 o/e = 0.75 (0.71 - 0.79) 0 → 1
pLoF	110.6	5	pLI = 1 o/e = 0.05 (0.02 - 0.1) 0 → 1

exome genome Metric: Mean Save plot

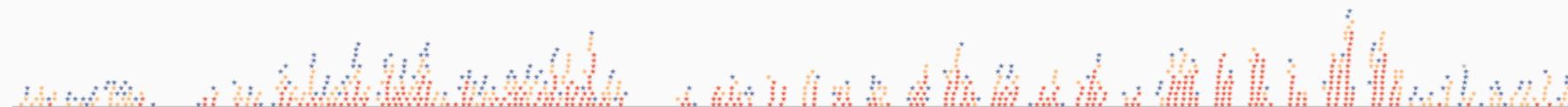


ClinVar variants

Pathogenic / likely pathogenic only Uncertain significance / conflicting only Benign / likely benign only Other only ?

Collapse to bins

ClinVar variants (812)



gnomAD variants

gnomAD v2.1.1 variants (2448)



Viewing in table

176,562,030 176,563,906 176,637,018 176,638,199 176,656,801 176,671,211 176,684,096 176,700,664 176,715,762 176,721,350 176,722,532

pLoF only Missense only Synonymous only Other only

Exomes Genomes SNVs Indels Filtered variants ?

Search variant table

Export variants to CSV

Note Only variants located in or within 75 base pairs of a coding exon are shown here. To see variants in UTRs or introns, use the region view.

The table below shows the HGVS consequence and VEP annotation for each variant's most severe consequence across all transcripts in this gene. Cases where the most severe consequence occurs in a non-canonical transcript are denoted with *. To see consequences in a specific transcript, use the transcript view.

Variant ID	Source	HGVS Consequence	VEP Annotation	Clinical Significance	Flags	Allele Count	Allele Number	Allele Frequency	Number of Heterozygote
5-176562031-AGAGTC-A	G	c.-17_51_17_47delGAGTC	● intron			2	31382	6.37e-5	0
5-176562058-C-A	G	c.-17_30C>A	● intron			1	31394	3.19e-5	0
5-176562062-T-G	E	c.-17_26T>G	● intron			2	250150	8.00e-6	0
5-176562065-A-G	F	c.-17_23A>G	● intron			3	250392	1.20e-5	0



Filter	Exomes	Genomes	Total
	Pass	Pass	
Allele Count	88	38	126
Allele Number	249800	31406	281206
Allele Frequency	0.0003523	0.001210	0.0004481
Popmax Filtering AF (95% confidence)	0.004016	0.003068	
Number of homozygotes	0	0	0

This variant is multiallelic. Other alt alleles are:

- 5-176562864-C-G

Annotations

This variant falls on 8 transcripts in 1 gene.

missense intron

- NSD1 • NSD1
- ENST00000439151 * • ENST00000347982
- HGVSp: p.Leu254Phe
- Polyphen: possibly_damaging
- SIFT: tolerated_low_confidence
- ENST00000361032 • ENST00000361032
- HGVSp: p.Leu254Phe
- Polyphen: possibly_damaging
- SIFT: tolerated_low_confidence

Population Frequencies

Population	Allele Count	Allele Number	Number of Homozygote
African	115	24888	0
Other	3	7168	0
Latino	7	35024	0
European (non-Finnish)	1	128728	0
Ashkenazi Jewish	0	10264	0
East Asian	0	19910	0
European (Finnish)	0	25070	0
South Asian	0	30154	0
Female	72	128682	0
Male	54	152524	0
Total	126	281206	0

References

- dbSNP (rs149334244)
- UCSC
- ClinVar (159441)

Report

- Report this variant
- Request additional information



Population Frequencies

Population	Allele Count	Allele Number	Number of Homozygotes	Allele Frequency
African	115	24888	0	0.004621
Other	3	7168	0	0.0004185
Latino	7	35024	0	0.0001999
Overall	1	128728	0	0.000007768
Other non-Finnish European	1	32984	0	0.00003032
Bulgarian	0	2666	0	0.000
Estonian	0	4836	0	0.000
North-western European	0	50588	0	0.000
Southern European	0	11578	0	0.000
Swedish	0	26076	0	0.000
Male	1	72028	0	0.00001388
Female	0	56700	0	0.000
Ashkenazi Jewish	0	10264	0	0.000
Overall	0	19910	0	0.000
Japanese	0	150	0	0.000
Korean	0	3818	0	0.000
Other East Asian	0	14382	0	0.000
Female	0	9840	0	0.000
Male	0	10070	0	0.000
European (Finnish)	0	25070	0	0.000
South Asian	0	30154	0	0.000
Female	72	128682	0	0.0005595
Male	54	152524	0	0.0003540
Total	126	281206	0	0.0004481

Additional features in gnomAD v3 (GRCh38)

In silico predictors

- ● REVEL: 0.325
- ● CADD: 21.9
- ● SpliceAI: 0.00 (No consequence)
- ● PrimateAI: 0.517

Released last week (also in gnomAD v2)

ClinVar

ClinVar Variation ID 96079

Conditions not specified, Beckwith-Wiedemann syndrome, History of neurodevelopmental disorder

Clinical significance Benign

Review status criteria provided, multiple submitters, no conflicts (2 stars)

Last evaluated 2019-12-31

[See all 5 submissions or find more information on the ClinVar website.](#)

Population Frequencies

gnomAD	HGDP	1KG			
Population	Allele Count	Allele Number	Number of Homozygotes	Allele Frequency	▼
▶ South Asian	262	4832	12	0.05422	
▶ Other	52	2092	0	0.02486	
▶ African/African-American	1027	41440	14	0.02478	
▶ European (non-Finnish)	1573	68026	21	0.02312	
▶ Middle Eastern	7	316	0	0.02215	
▶ Latino/Admixed American	326	15290	3	0.02132	
▶ East Asian	89	5206	3	0.01710	
▶ European (Finnish)	124	10610	0	0.01169	
▶ Amish	10	912	0	0.01096	
▶ Ashkenazi Jewish	14	3470	0	0.004035	
XX	1796	77830	26	0.02308	
XY	1688	74364	27	0.02270	
Total	3484	152194	53	0.02289	

Additional features in gnomAD v3 (GRCh38)

- Two open access datasets
- Human Genome Diversity Project (HGDP)
- 1000Genomes (1KG)

Population	Allele Count	Allele Number	Number of Homozygotes	Allele Frequency
Central/South Asian	25	332	1	0.07530
African	6	96	0	0.06250
Middle Eastern	10	308	0	0.03247
East Asian	14	456	0	0.03070
European	5	298	0	0.01678
Native American	0	70	0	0.000
XX	16	568	0	0.02817
XY	44	992	1	0.04435
Total	60	1560	1	0.03846

Warning Because of low sample sizes for HGDP populations, allele frequencies may not be representative.

Additional features in gnomAD v3 (GRCh38)

Two open access datasets
Human Genome Diversity
Project (HGDP)

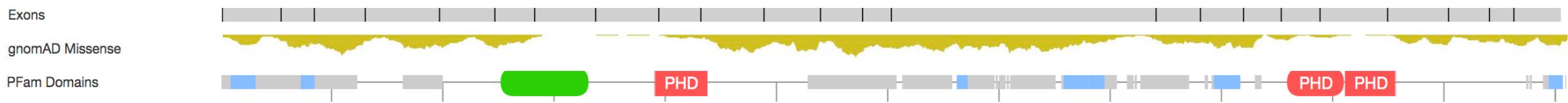
1000Genomes (1KG)

Population	Allele Count	Allele Number	Number of Homozygotes	Allele Frequency
Central/South Asian	25	332	1	0.07530
Overall	6	96	0	0.06250
Bantu (Kenya)	3	14	0	0.2143
Mandenka	2	42	0	0.04762
Yoruba	1	38	0	0.02632
African				
Bantu (South Africa)	0	2	0	0.000
XX	2	32	0	0.06250
XY	4	64	0	0.06250
Middle Eastern	10	308	0	0.03247
East Asian	14	456	0	0.03070
European	5	298	0	0.01678
Native American	0	70	0	0.000
XX	16	568	0	0.02817
XY	44	992	1	0.04435
Total	60	1560	1	0.03846

Warning Because of low sample sizes for HGDP populations, allele frequencies may not be representative.

Constraint

- Where do we see depletion of variation across human populations?



Conservation is between species
Constraint is within the human species

Developed a mutational model to calculate expected number of variants per gene

Category	<u>Exp. SNVs</u>	<u>Obs. SNVs</u>
Synonymous	<u>525.3</u>	523
Missense	<u>1427.7</u>	1065
pLoF	<u>110.6</u>	5



Kaitlin Samocha



Konrad Karczewski

Daniel MacArthur
Mark Daly

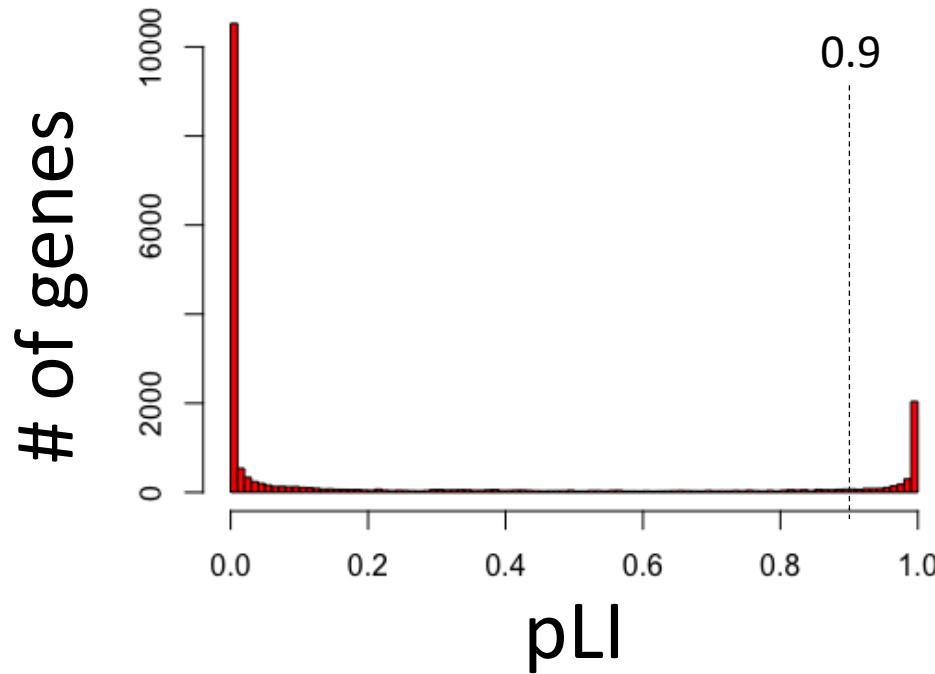
Constraint metrics

Category	Exp. SNVs	Obs. SNVs	Constraint metrics
Synonymous	525.3	523	$Z = 0.08$ $\text{o/e} = 1 (0.93 - 1.07)$
Missense	1427.7	1065	$Z = 3.41$ $\text{o/e} = 0.75 (0.71 - 0.79)$
pLoF	110.6	5	$\text{pLI} = 1$ $\text{o/e} = 0.05 (0.02 - 0.1)$ ← LOEUF

- Used z-scores to assess depletion from expectation
- For pLoF variants, correlation with gene length (short genes under powered)

Constraint metrics

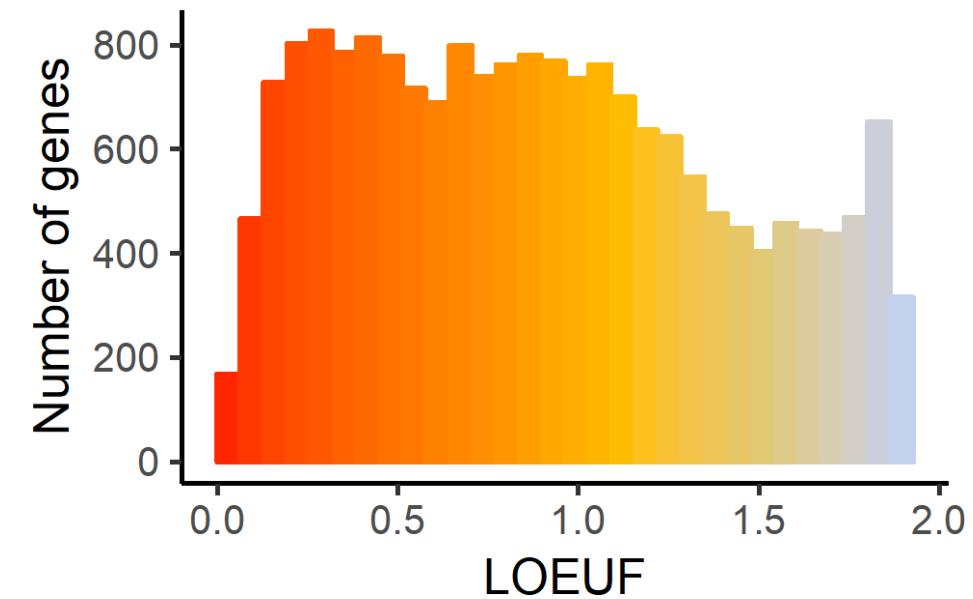
pLI = Probability of loss function intolerance
Dichotomous metric (>0.9 LoF constrained)



Samocha *et al.* 2014; Lek *et al.* 2016

Category	Exp. SNVs	Obs. SNVs	Constraint metrics
Synonymous	525.3	523	Z = <u>0.08</u> o/e = <u>1</u> (0.93 - 1.07)
Missense	1427.7	1065	Z = 3.41 o/e = <u>0.75</u> (0.71 - 0.79)
pLoF	110.6	5	pLI = <u>1</u> o/e = <u>0.05</u> (0.02 - 0.1)

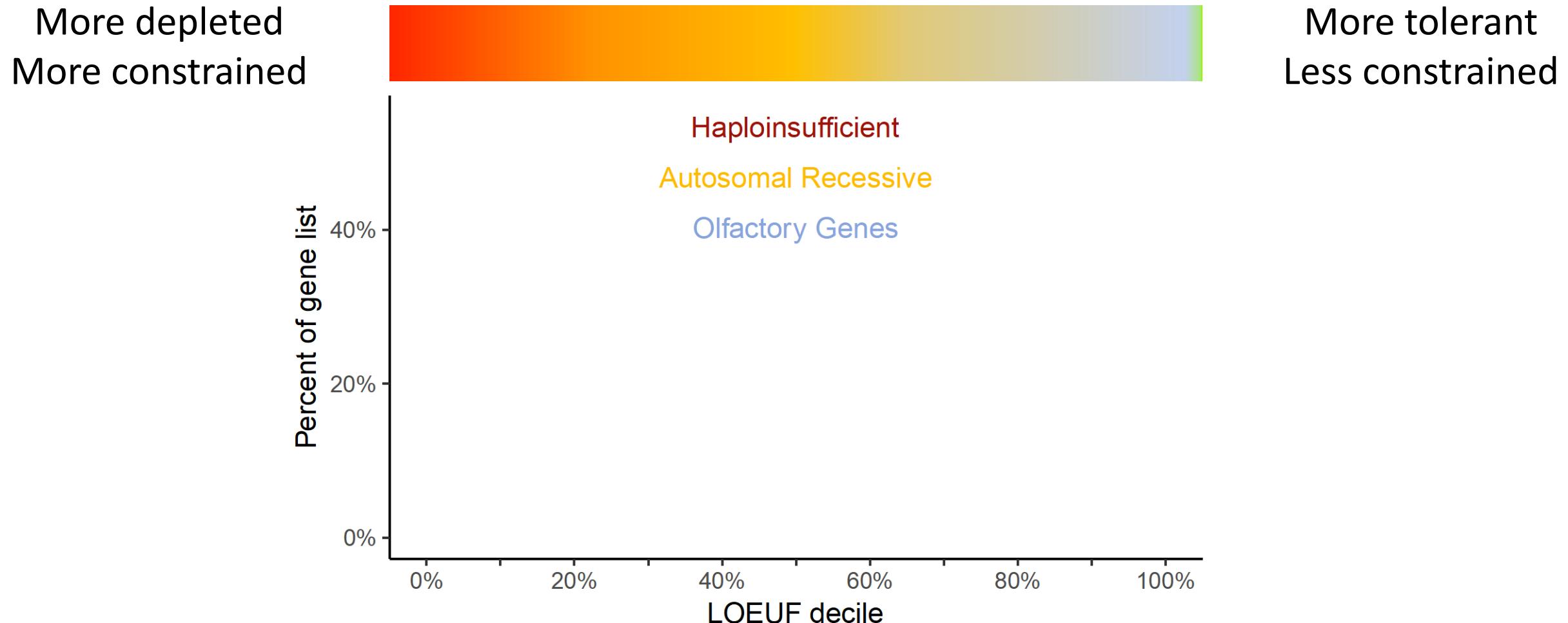
LOEUF = Loss-of-function observed/expected upper bound fraction
Continuous metric (lower scores more LoF constrained)



Karczewski *et al.*, Nature, 2020

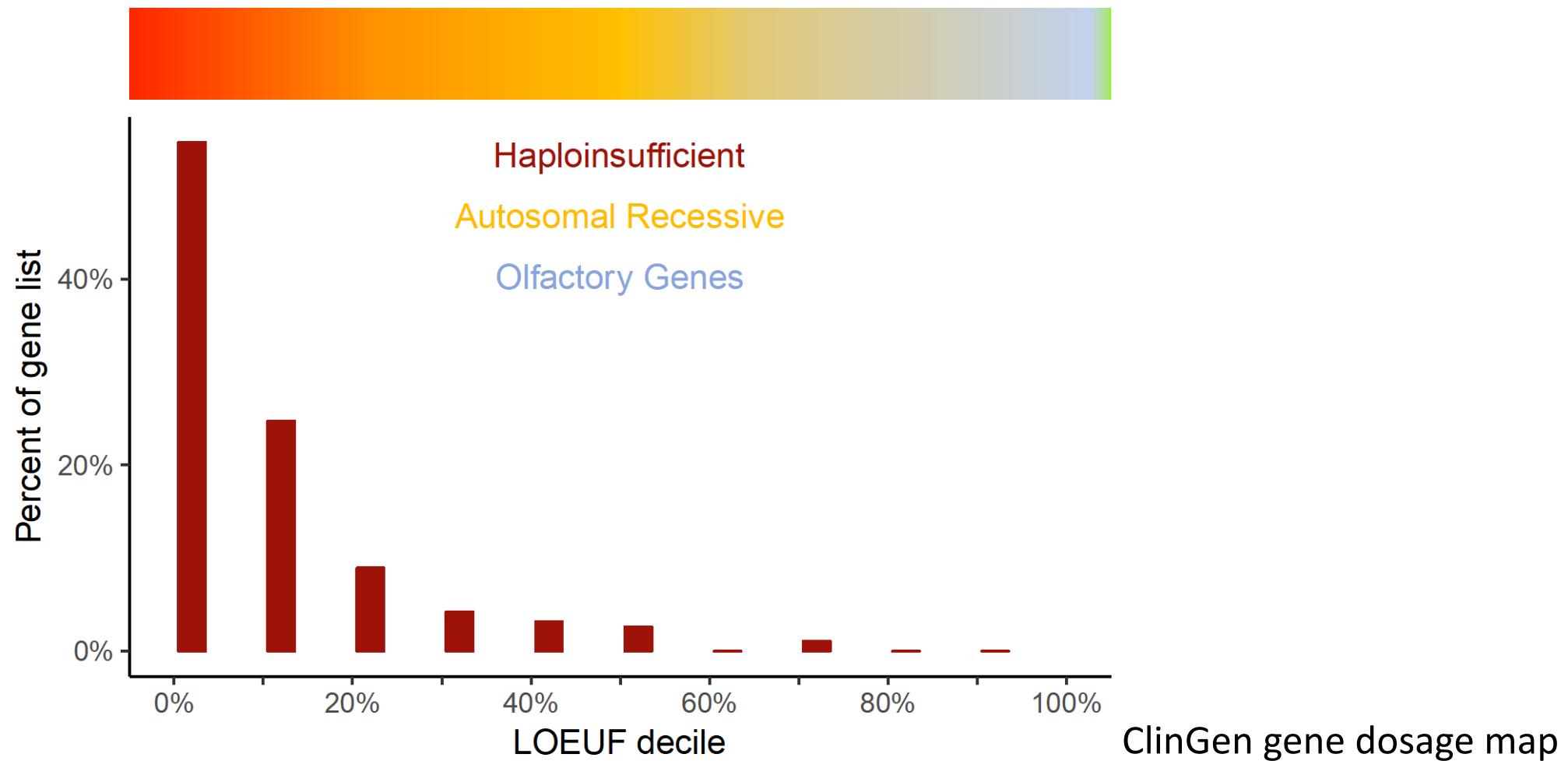
Resolving the spectrum of LoF intolerance

- Binning this spectrum into deciles



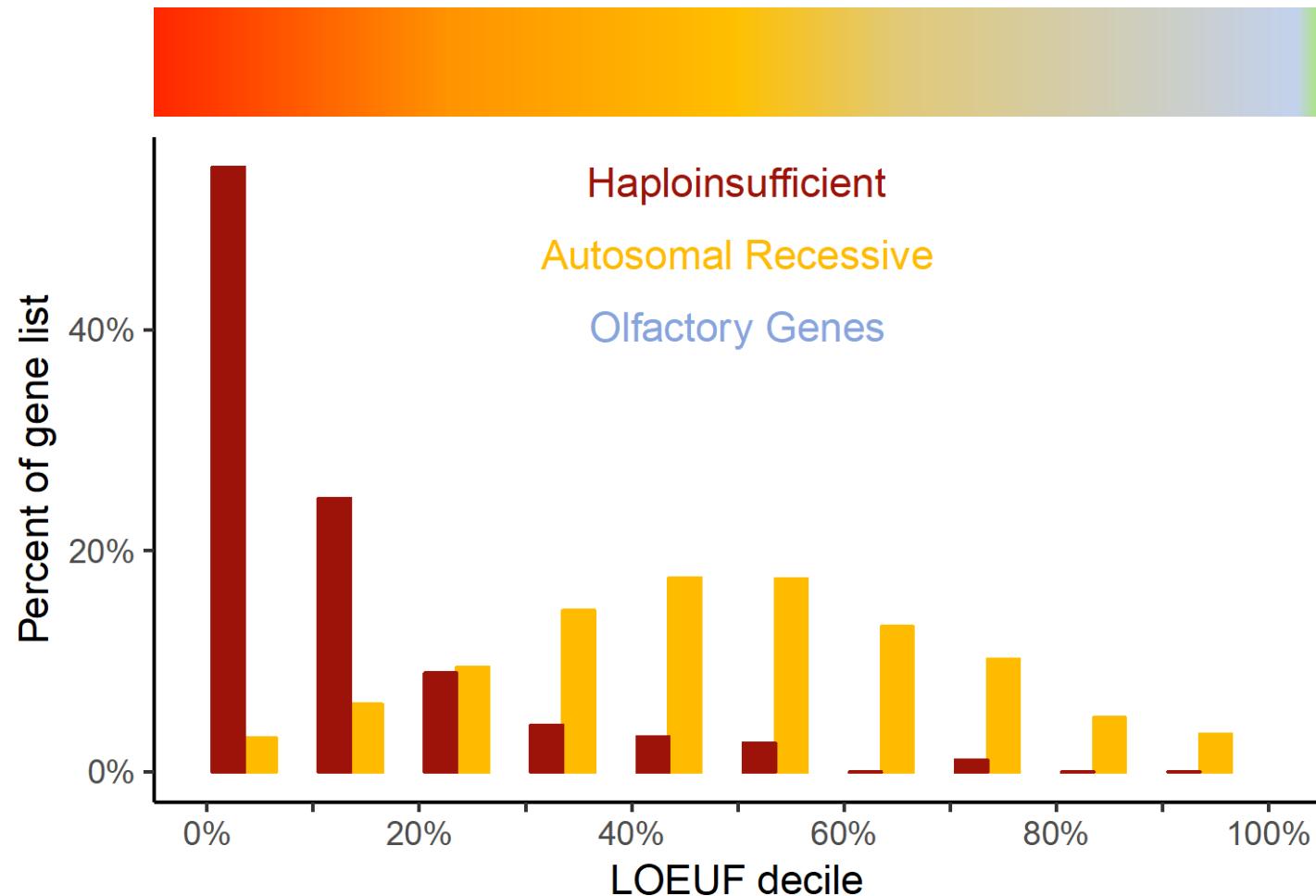
Resolving the spectrum of LoF intolerance

- Known haploinsufficient genes where loss of one copy results in disease



Resolving the spectrum of LoF intolerance

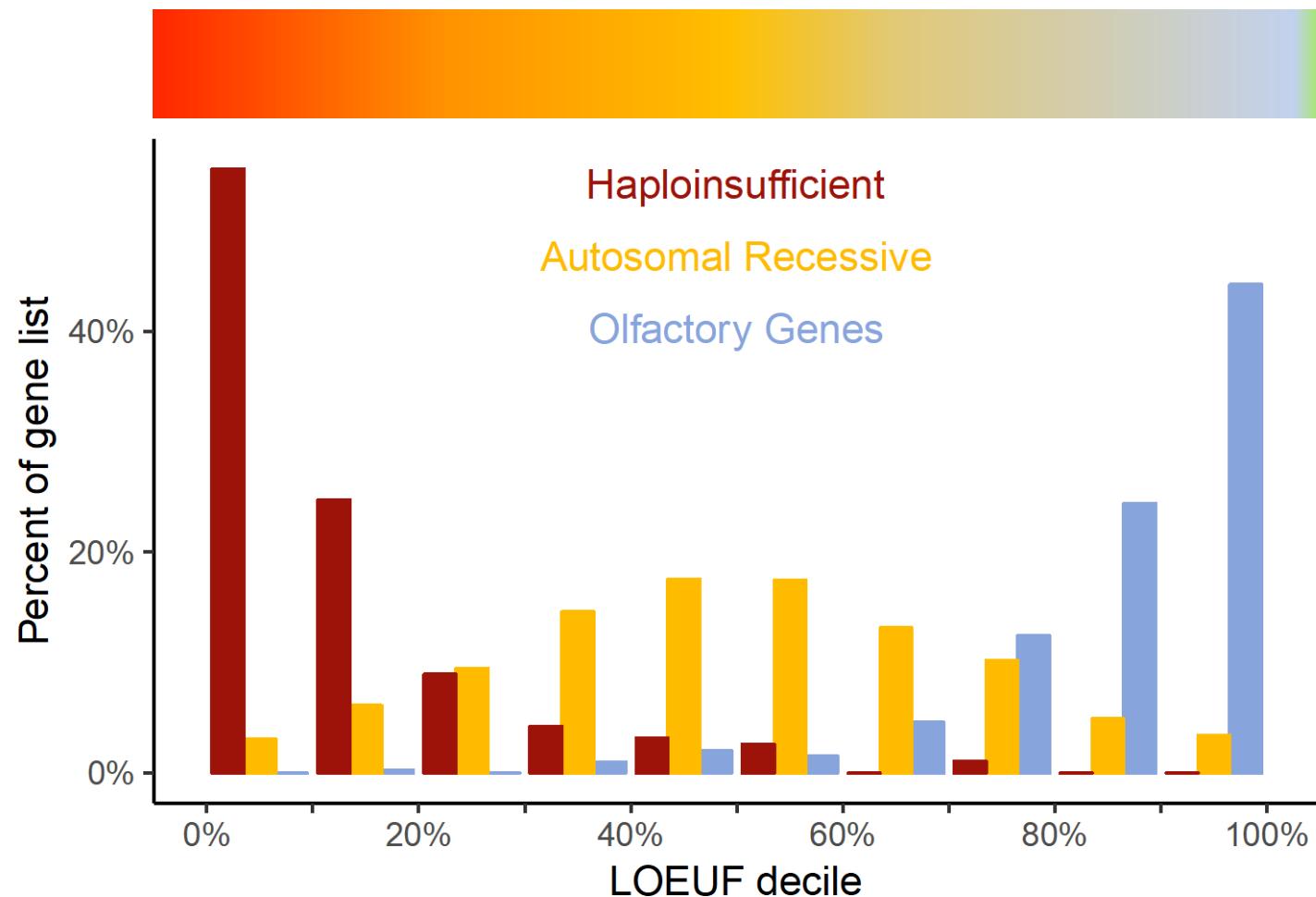
- Autosomal recessive genes are centered around 60% of expected



Gene list from:
Blekhan et al., 2008
Berg et al., 2013

Resolving the spectrum of LoF intolerance

- Some genes, e.g. olfactory receptors, are unconstrained



Important notes on constraint

- **Primarily for dominant disease genes;** do not expect to see strong constraint signal (depletion of pLoF variation) for recessive disease genes
- Not necessarily a sign of a strong phenotype – can see evidence of constraint from mild degrees of negative selection
- Selection occurs through reproduction so deleterious effect must be before/during reproductive years
 - Example: *BRCA1* does not show evidence of loss of function constraint

BRCA1
Autosomal dominant
Breast and ovarian cancer

Category	Exp. SNVs	Obs. SNVs	Constraint metrics	
Synonymous	341	329	Z = 0.51 o/e = 0.96 (0.88 - 1.06)	
Missense	941.3	891	Z = 0.58 o/e = 0.95 (0.9 - 1)	
pLoF	75.2	55	pLI = 0 o/e = 0.73 (0.59 - 0.92)	

Regional missense constraint (Samocha et al.)

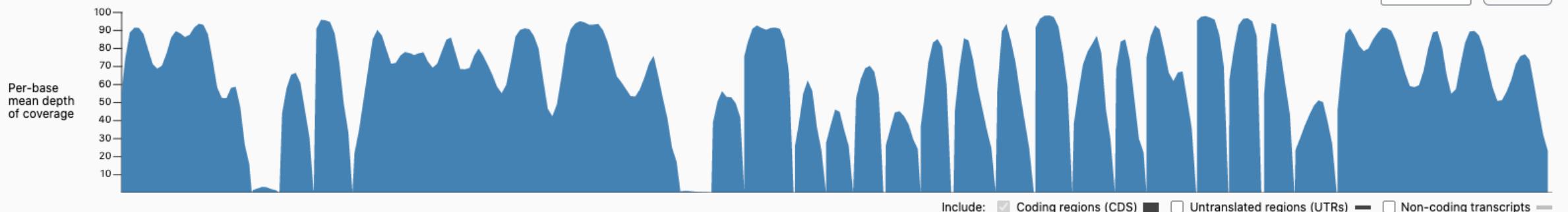
<https://www.biorxiv.org/content/early/2017/06/12/148353>

[more](#)

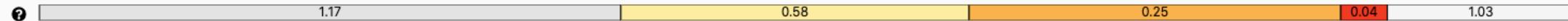
Constraint ρ

Category	Exp. SNVs	Obs. SNVs	Constraint metric
Synonymous	323.8	345	Z = -0.73
Missense	833.7	693	Z = 2.38
pLoF	74.8	4	pLI = 1

exome genome Metric: Mean Save plot

[Show transcripts](#)[Show tissues](#)

Regional missense constraint



ClinVar variants

Pathogenic / likely pathogenic only Uncertain significance / conflicting only Benign / likely benign only Other only

[Collapse to bins](#)

ClinVar variants (812)

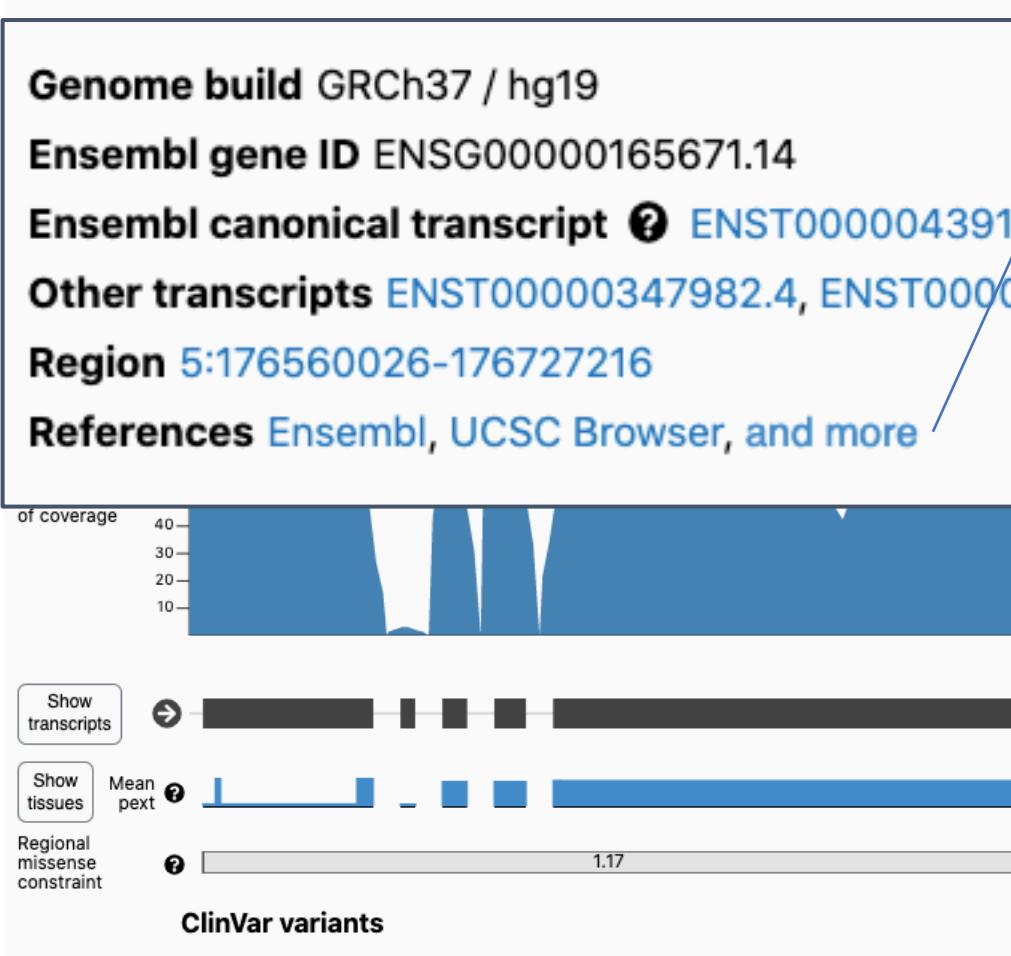


gnomAD variants

ExAC v1.0 variants







References for NSD1

- Ensembl
- UCSC Browser
- GeneCards
- OMIM
- DECIPHER
- ClinGen
- Gene Curation Coalition (GenCC)
- HGNC

Ok

DECIPHER Protein View: *NSD1*

<https://decipher.sanger.ac.uk/gene/NSD1/overview/protein-info>

NSD1: Q96L73 2696aa

?

Links ▾

Settings



Expression data from GTEx in gnomAD

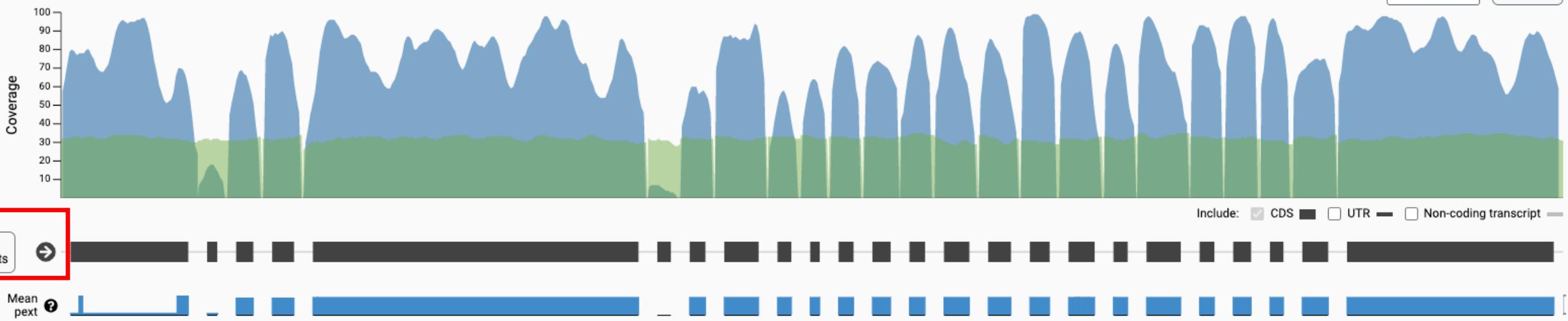
NSD1 nuclear receptor binding SET domain protein 1

Dataset gnomAD v2.1.1 ▾ gnomAD SVs v2.1 ▾ ?

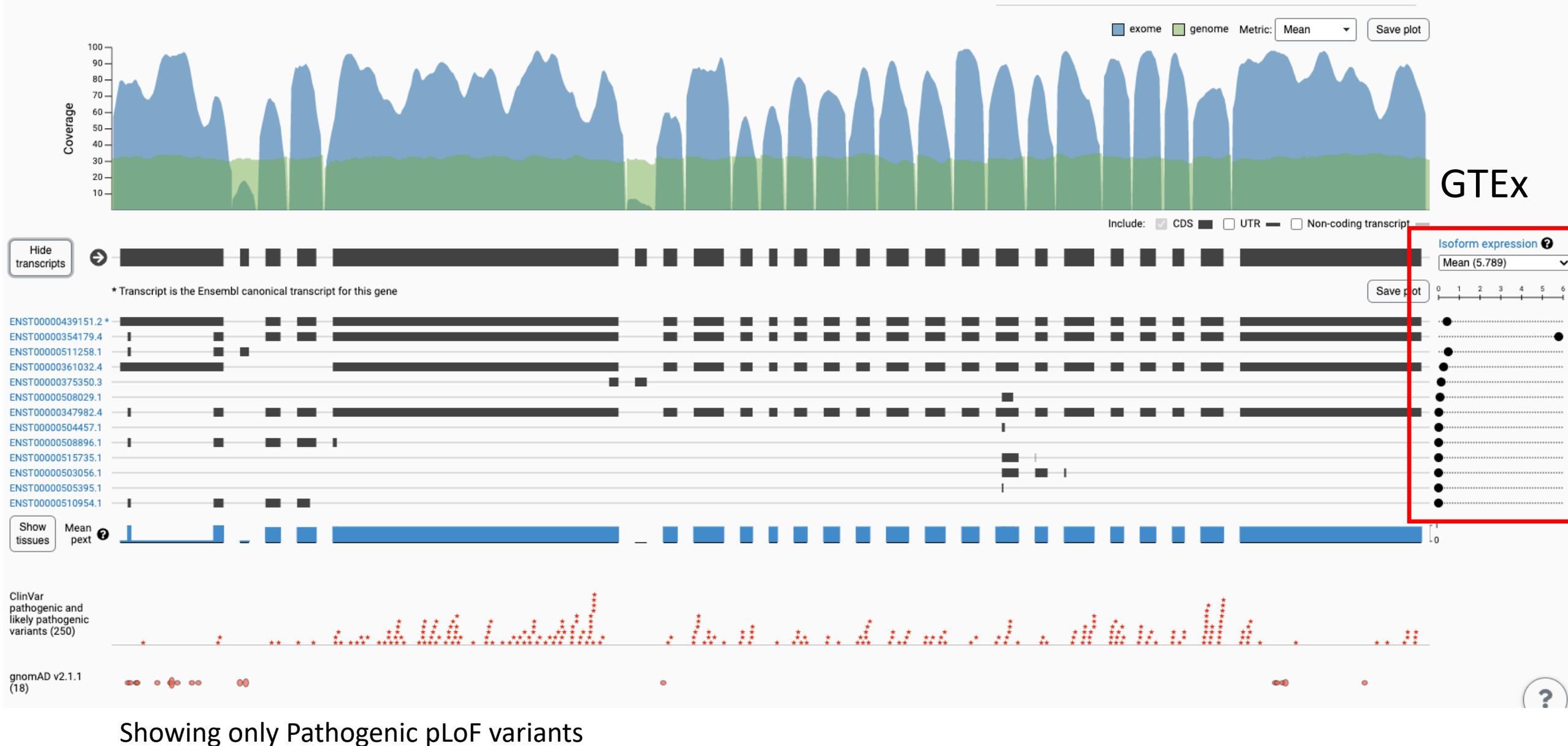
Genome build GRCh37 / hg19
Ensembl gene ID ENSG00000165671.14
Ensembl canonical transcript ⓘ ENST00000439151.2
Region 5:176560026-176727216
References Ensembl, UCSC Browser, and more

Constraint ⓘ

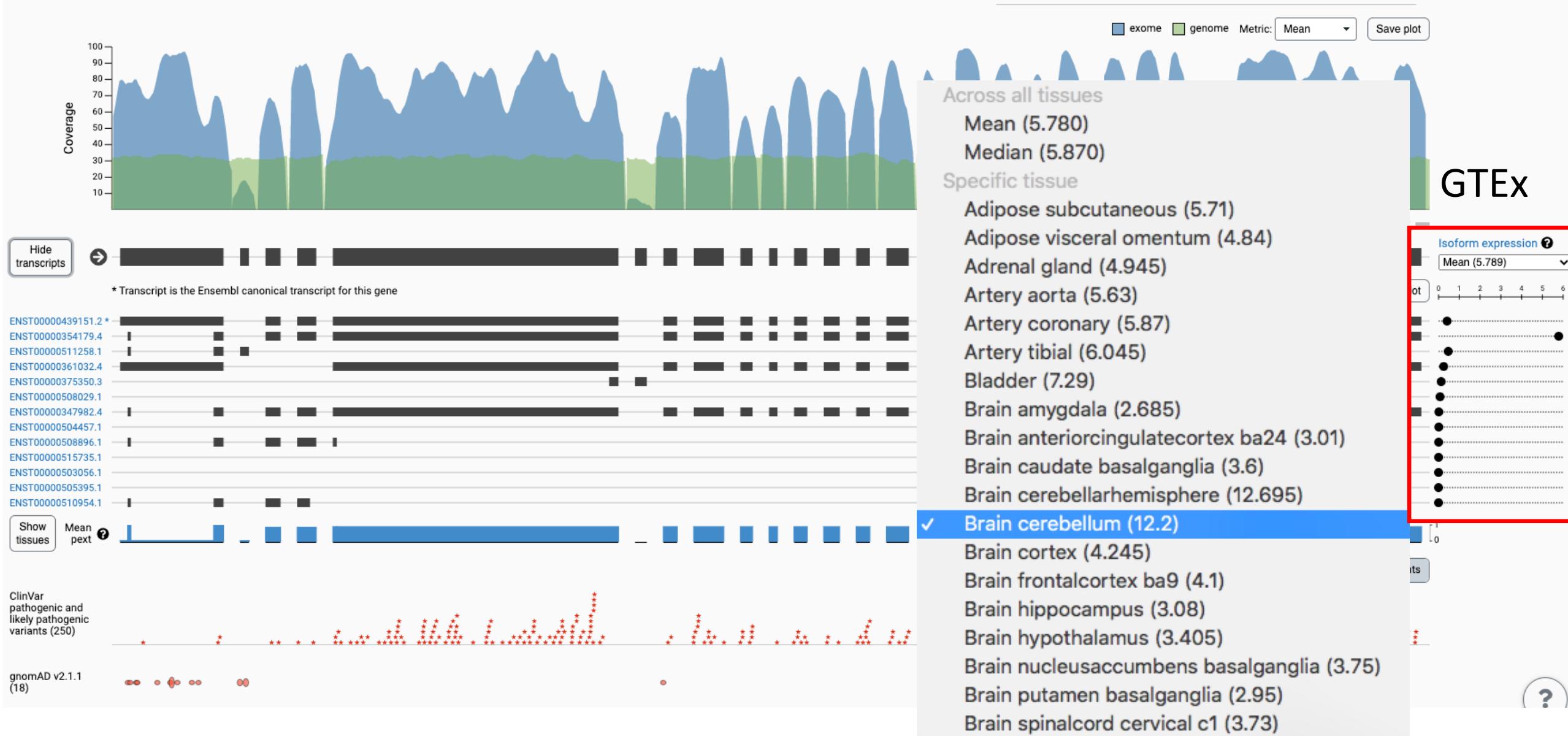
Category	Exp. SNVs	Obs. SNVs	Constraint metrics
Synonymous	525.3	523	Z = 0.08 o/e = 1 (0.93 - 1.07) 0 ⚡ 1
Missense	1427.7	1065	Z = 3.41 o/e = 0.75 (0.71 - 0.79) 0 ⚡ 1
pLoF	110.6	5	pLI = 1 o/e = 0.05 (0.02 - 0.1) 0 ⚡ 1



NSD1



NSD1



Home Datasets Expression QTLs & Browsers Sample Data Documentation

NSD1

2018-05-10 Access GTEx Biospecimens Here We have updated our biospecimen search page. New features include a shopping cart for streamlined sample requests: <https://gtexportal.org/home/samplesPage>

Resource Overview

Current Release (V8)

- Tissue & Sample Statistics
- Tissue Sampling Info (Anatomogram)
- Access & Download Data
- Release History
- How to cite GTEx?

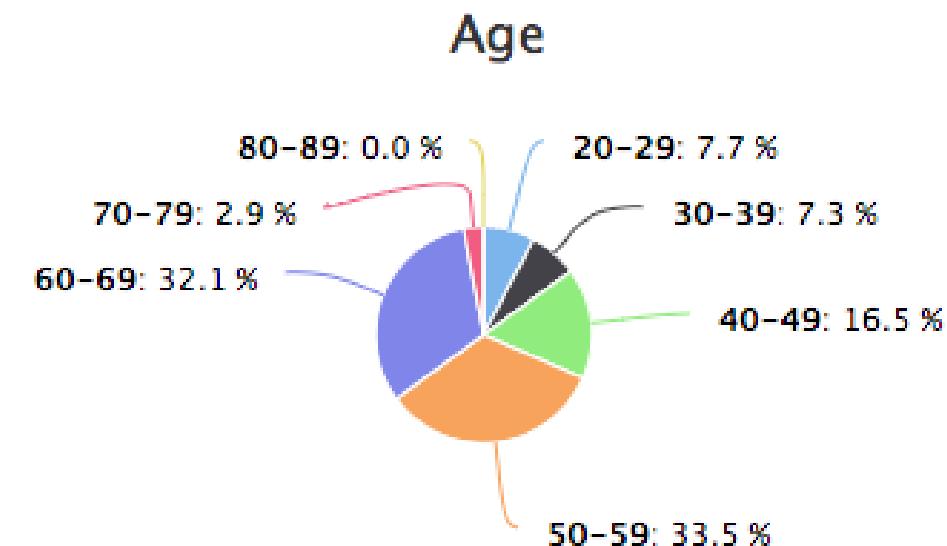
The Genotype-Tissue Expression (GTEx) project is an ongoing effort to build a comprehensive public resource to study tissue-specific gene expression and regulation. Samples were collected from 54 non-diseased tissue sites across nearly 1000 individuals, primarily for molecular assays including WGS, WES, and RNA-Seq. Remaining samples are available from the GTEx Biobank. The GTEx Portal provides open access to data including gene expression, QTLs, and histology images.

Explore GTEx

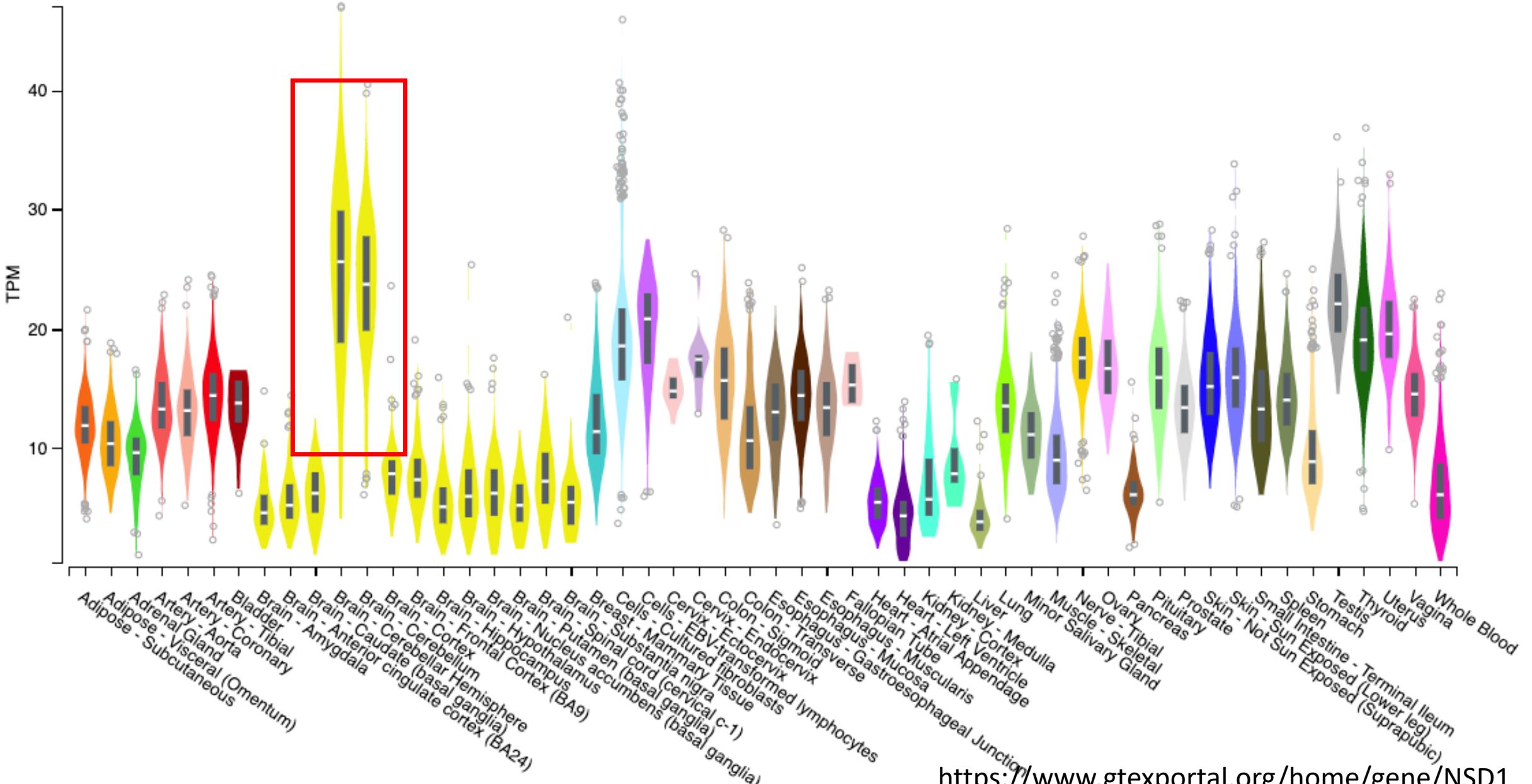
Browse		By gene ID	Browse and search all data by gene
		By variant or rs ID	Browse and search all data by variant
		By Tissue	Browse and search all data by tissue
		Histology Image Viewer	Browse and search GTEx histology images
Expression		Multi-Gene Query	Browse and search expression by gene and t
		Top 50 Expressed Genes	Visualize the top 50 expressed genes in each
		Transcript Browser	Visualize transcript expression and isoform st

V7 Release	# Tissues	# Donors	# Samples
Total	53	714	11688

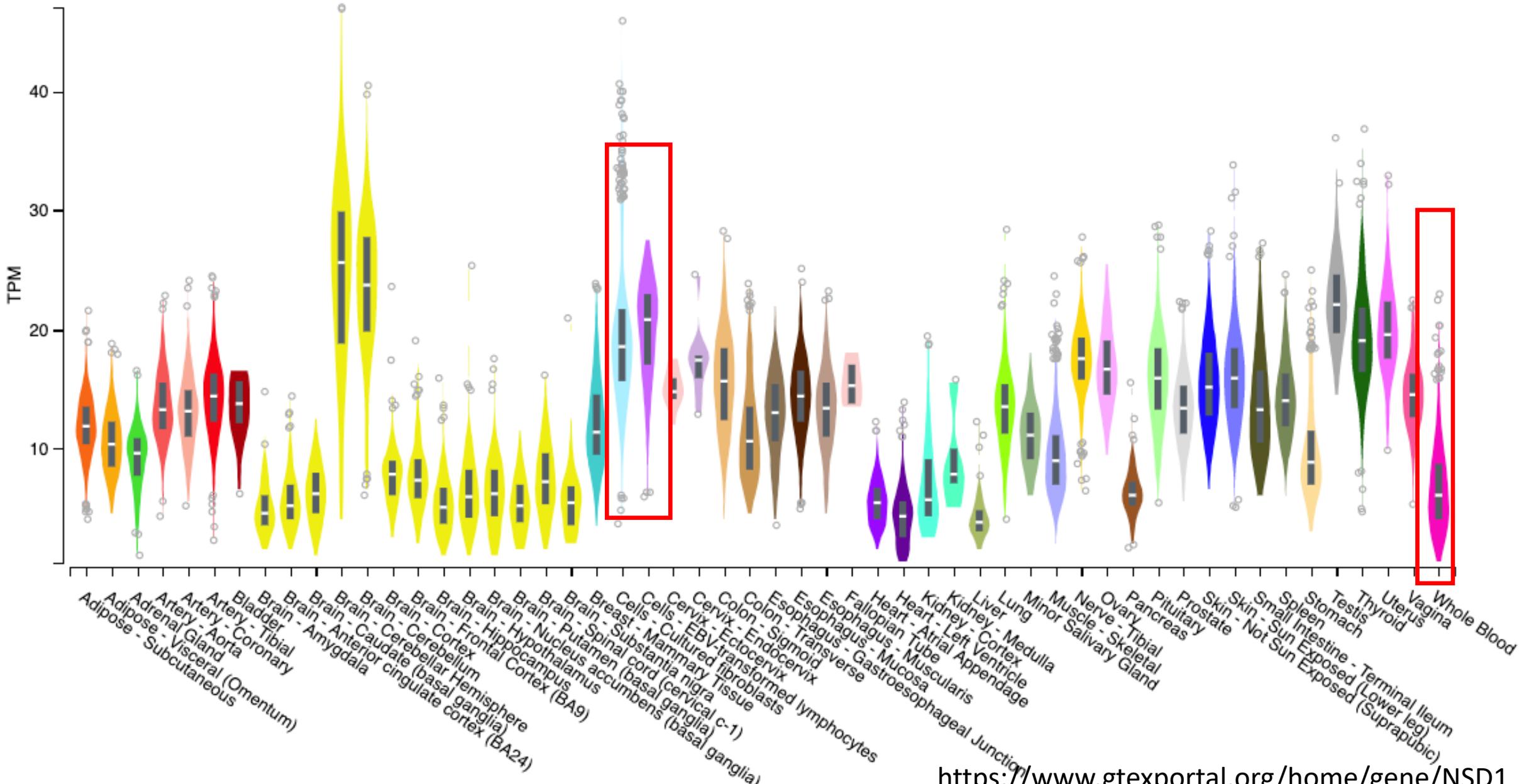
Gene expression Transcript expression eQTLs



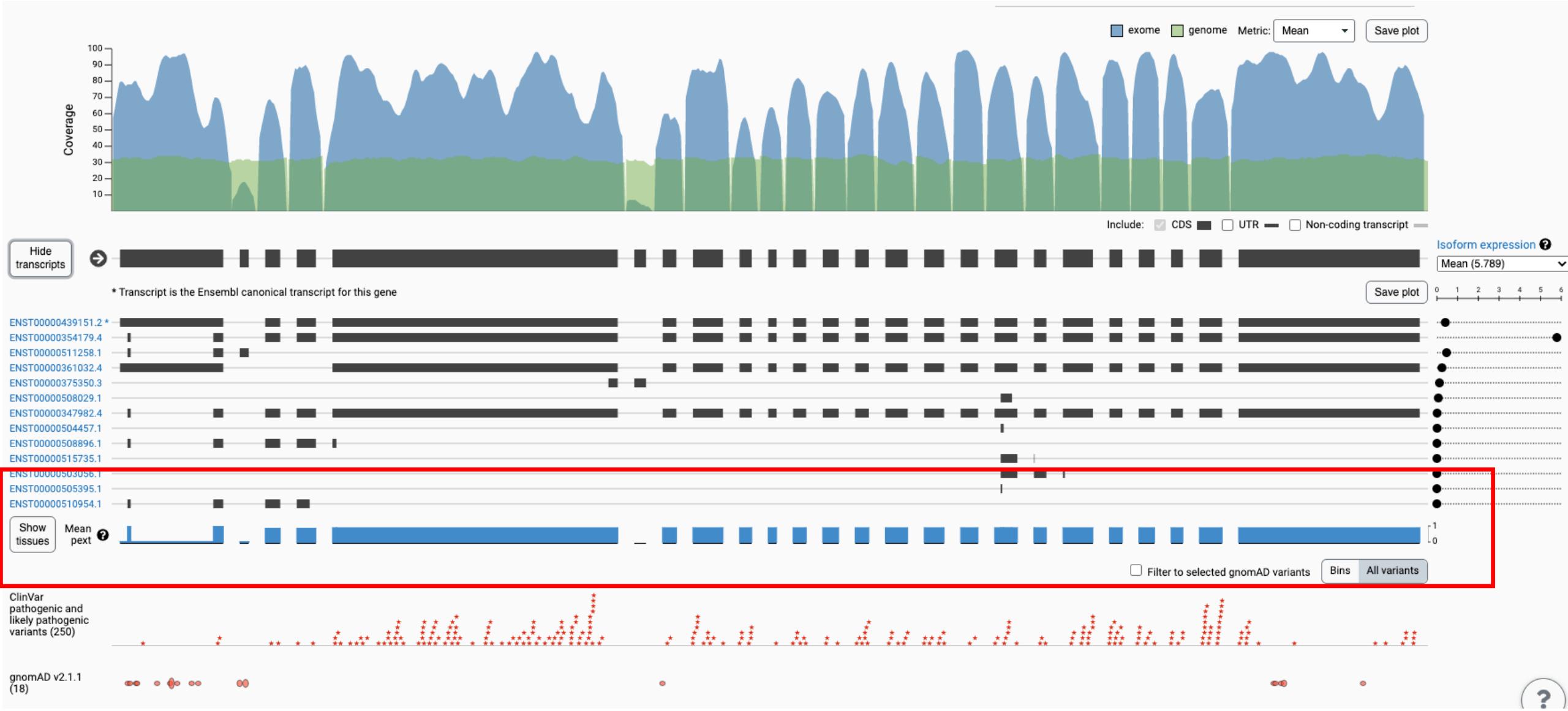
Example gene expression across tissues: *NSD1*



Example gene expression across tissues: *NSD1*



Transcript aware expression analysis - pext score



Exploring subsets

Single nucleotide variant: 5-176631141-C-T (

	Exomes	Genomes	Total
Filter	Pass	No variant	
Allele Count	1		1
Allele Number	251484		251484
Allele Frequency	0.000003976		0.000003976
Popmax Filtering AF ? (95% confidence)	—		
Number of homozygotes	0		0

gnomAD v2.1.1 ?

gnomAD v2.1.1
141,456 samples

gnomAD v2.1.1 (non-TOPMed)
135,743 samples

gnomAD v2.1.1 (non-cancer)
134,187 samples

gnomAD v2.1.1 (non-neuro)
114,704 samples

gnomAD v2.1.1 (controls)
60,146 samples

ExAC v1.0
60,706 samples

missense

- NSD1
- ENST00000439151.2

Ensembl canonical transcript for |
HGVSp: p.Arg362Trp

Population Frequencies ?

Population	Allele Count	Allele Number	Number of Homozygotes	Allele Frequency
▶ African	1	16256	0	0.00006152
▶ Latino	0	34590	0	0.000
▶ Ashkenazi Jewish	0	10080	0	0.000
▶ East Asian	0	18394	0	0.000
▶ European (Finnish)	0	21648	0	0.000

Exploring subsets

The screenshot shows the gnomAD browser interface. At the top left is the "gnomAD browser" logo. In the center is a search bar with "gnomAD v2.1.1" and a "Search" button. To the right of the search bar are links for "About", "News", "Downloads", "Terms", "Publications", and "Contact". A blue banner across the middle says "gnomAD v3.1 released!". Below the banner, the main content area has a title "Single nucleotide variant: 5-176631141-C-T (GRCh37)" and a dataset selector "Dataset: gnomAD v2.1.1 (non-cancer) ?". A message "Variant not found in selected subset." is displayed in the center of the page.

Cancer means germline sample (typically blood) for individual recruited to a study for a non-blood cancer

NSD1 nuclear receptor binding SET domain protein 1

Genome build GRCh37 / hg19

Ensembl gene ID ENSG00000165671

Canonical transcript ID ENST00000439151

Region 5:176560027-176727217

References Ensembl, UCSC Browser, and more

Gene page region view

5-176560027-176727217

Change

Dataset gnomAD v2.1.1

gnomAD SVs v2.1



Genome build GRCh37 / hg19

Region size 167,191 BP

References UCSC Browser

Zoom in

1.5x

3x

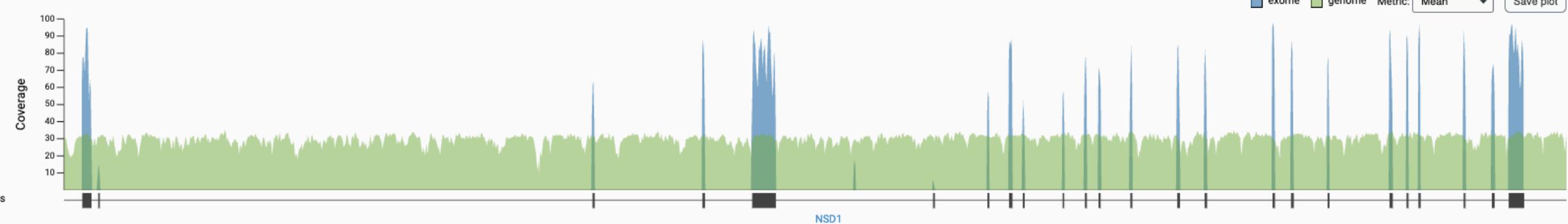
10x

Zoom out

1.5x

3x

10x



ClinVar
pathogenic
and likely
pathogenic
variants (335)
61-
0

Filter to selected gnomAD variants

Bins All variants

gnomAD v2.1.1
(15548)

A structural variation reference for medical and population genetics

Ryan L. Collins, Harrison Brand, [...] Michael E. Talkowski 

Nature 581, 444–451(2020) | Cite this article

29k Accesses | 38 Citations | 168 Altmetric | Metrics



Ryan
Collins



Harrison
Brand



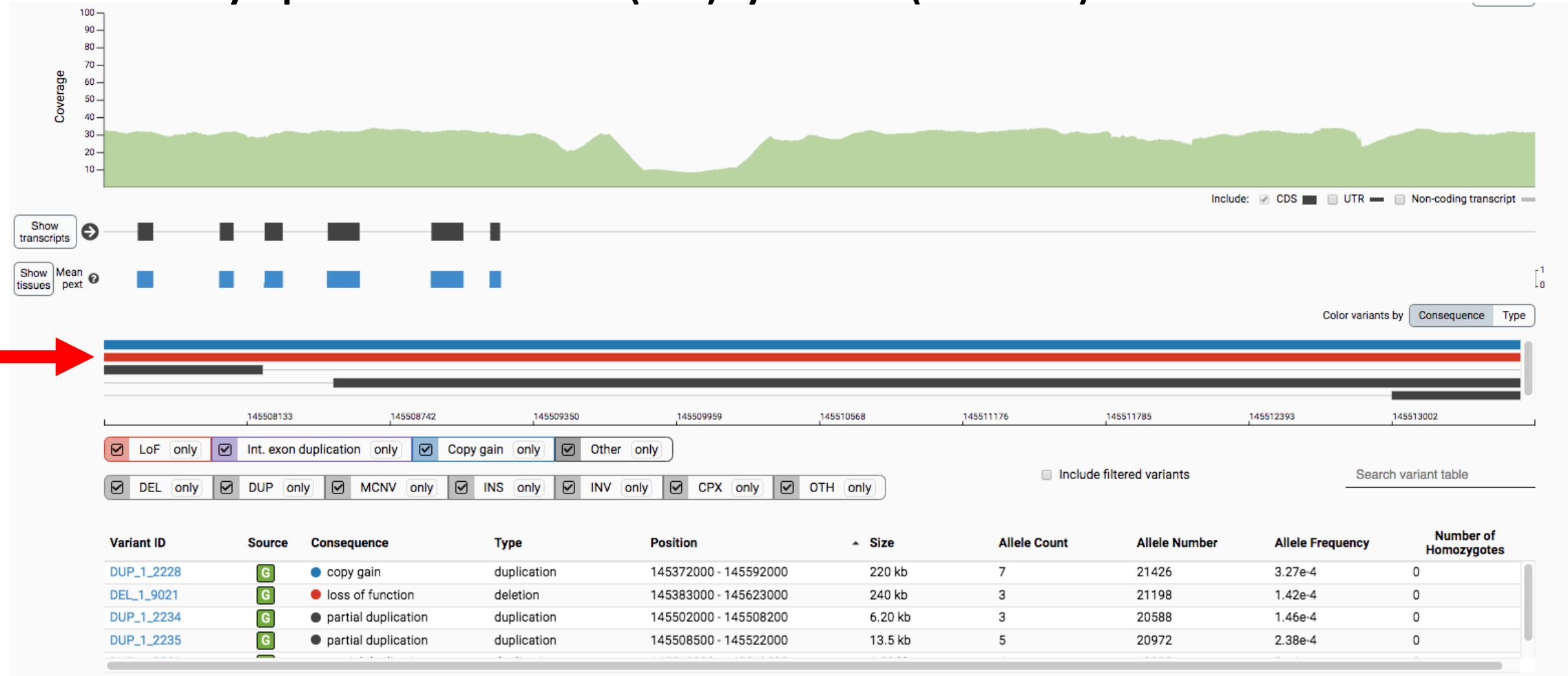
Mike
Talkowski

- gnomAD-SV callset from 10,847 whole genomes
- Multiple algorithms followed by integration step
- 433,371 unique structural variants with frequency data

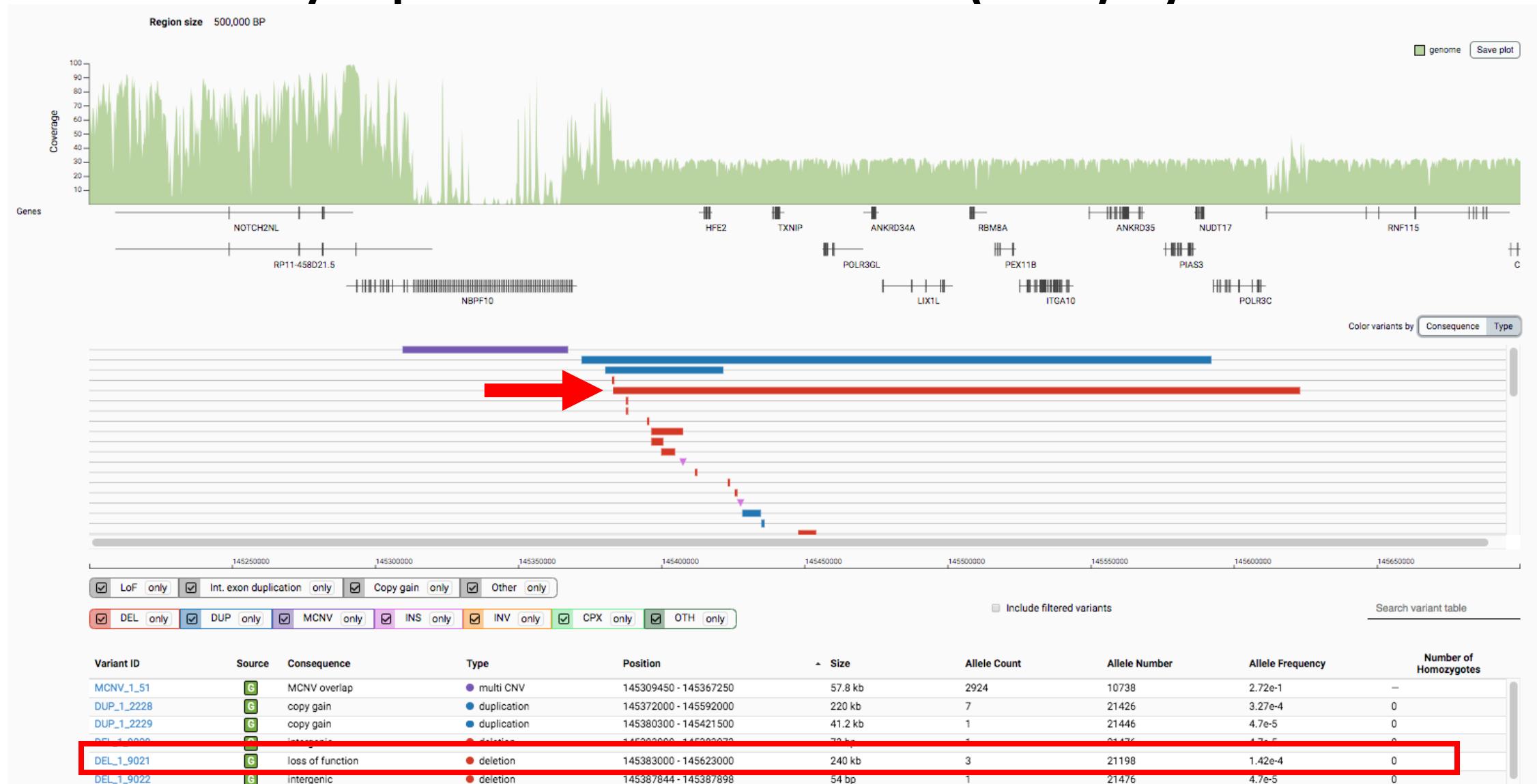
gnomAD Structural Variants: *RBM8A*

Dataset gnomAD v2.1.1 ▾ gnomAD SVs v2.1 ▾

Thrombocytopenia absent radius (TAR) syndrome (recessive)



Previously reported “common” deletion in Thrombocytopenia absent radius (TAR) syndrome



Structural variant: DEL_1_7663

Dataset

gnomAD SVs v2.1

Filter	Pass
Allele Count	3
Allele Number	21314
Allele Frequency	0.0001408
Quality score	999
Position	1:145383000-145833000
Size	450,000 bp
Class	deletion ?
Evidence	Read depth
Algorithms	Depth

Consequences

This variant has consequences in 16 genes.

loss of function [?](#)

- [ANKRD34A](#)
- [ANKRD35](#)
- [CD160](#)
- [and 13 more](#)

References

- [UCSC](#)

Report

- [Report this variant](#)

gnomAD SVs v2.1

10,847 samples

gnomAD SVs v2.1 (controls)

5,192 samples

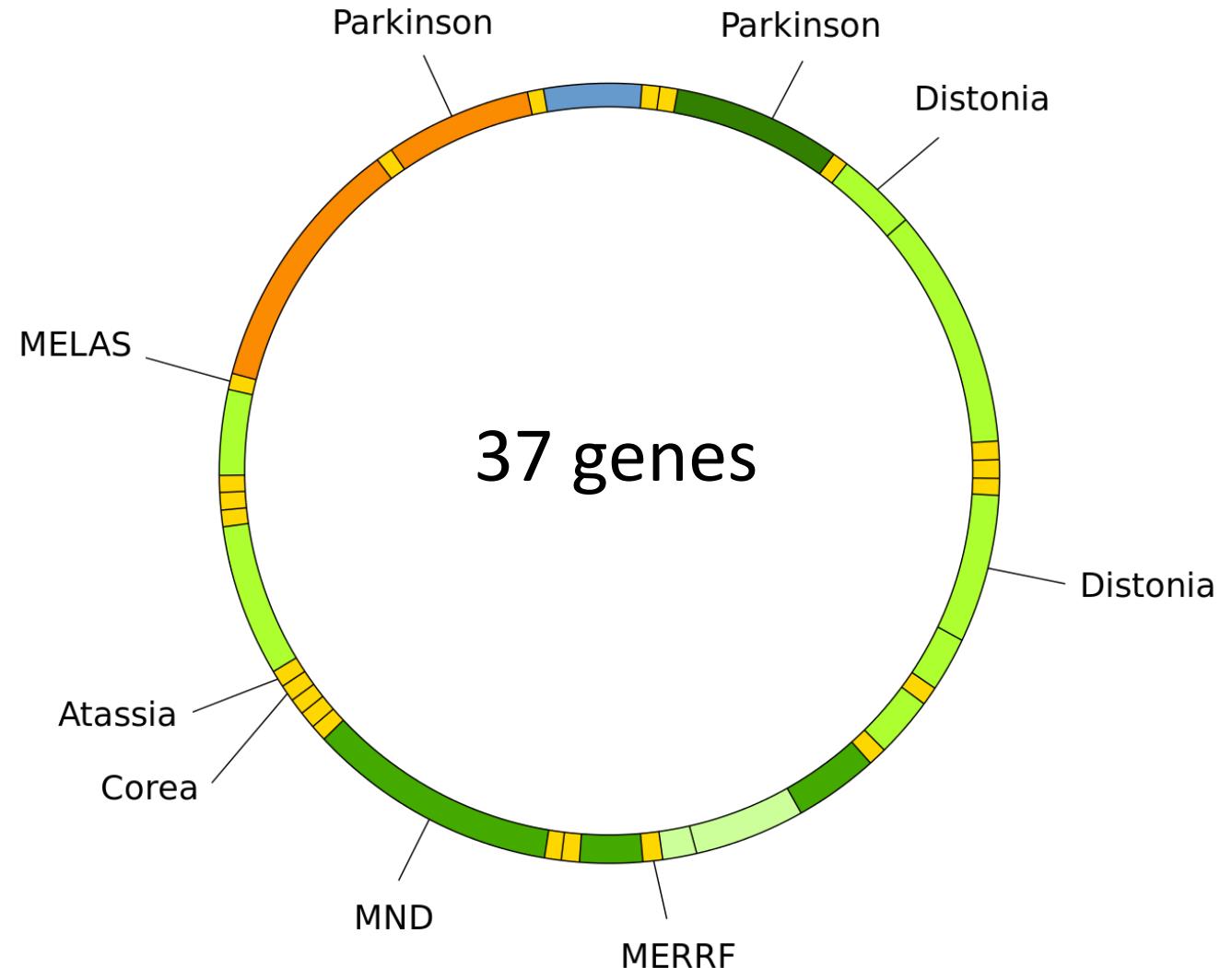
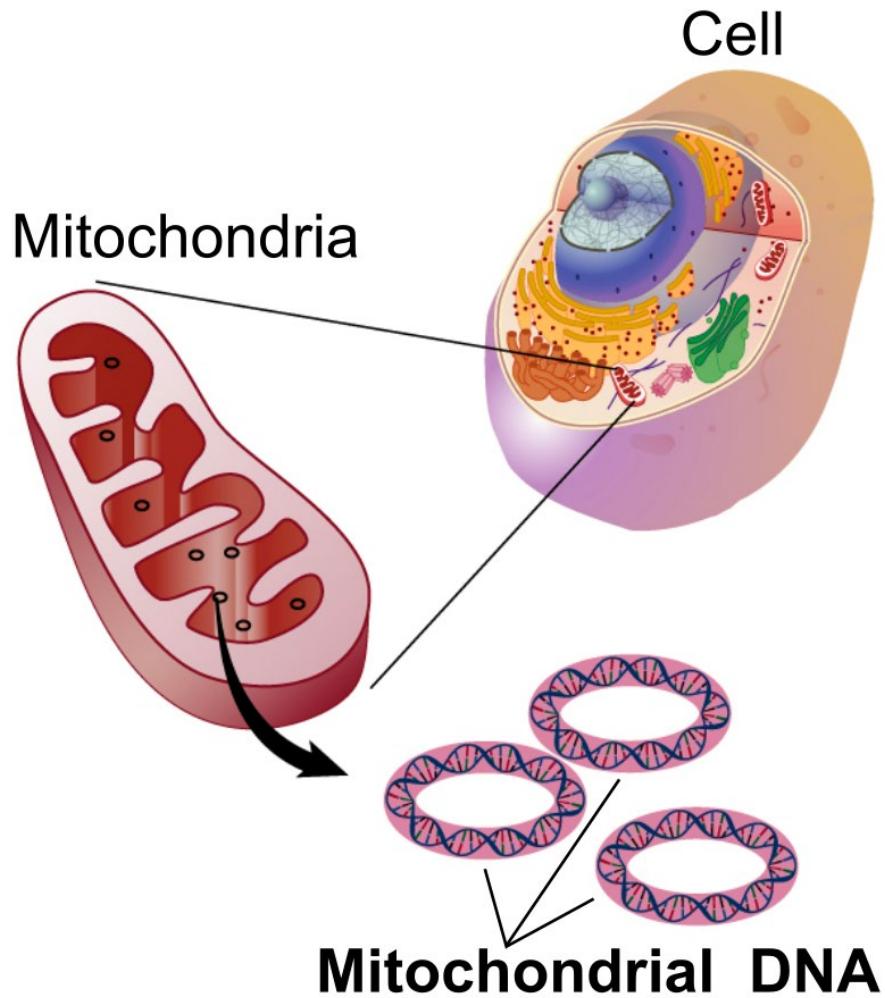
gnomAD SVs v2.1 (non-neuro)

8,342 samples

Population Frequencies

Population	Allele Count	Allele Number	Number of Homozygotes	Allele Frequency
European	2	7534	0	0.0002655
African	1	9406	0	0.0001063
Latino	0	1928	0	0.000
East Asian	0	2258	0	0.000
Other	0	188	0	0.000
Female	2	10416	0	0.0001920
Male	1	10852	0	0.00009215
Total	3	21314	0	0.0001408

Mitochondrial variants in gnomAD v3



Genome build GRCh38 / hg38

Ensembl gene ID ENSG00000198786.2

MANE Select transcript Not available

Ensembl canonical transcript Not available

Other transcripts ENST00000361567.2

Region M:12337-14148

References Ensembl, UCSC Browser, and more

Constraint

Constraint not yet available for gnomAD v3.

**ClinVar variants**
 Pathogenic / likely pathogenic only Uncertain significance / conflicting only Benign / likely benign only Other only

ClinVar variants (300)

**gnomAD variants**

gnomAD variants (1136)



Variant ID	Source	HGVS Consequence	VEP Annotation	Clinical Significance	Flags	Allele Number	Homoplasmic Allele Count	Homoplasmic Allele Frequency	Heteroplasmic Allele Count	Heteroplasmic Allele Frequency	Max observed heteroplasmy
M-12372-G-A		p.Leu12Leu	● synonymous	Likely pathogenic	Common Low Heteroplasmy	56353	8901	1.58e-1	8	1.42e-4	1
M-12372-G-C		p.Leu12Leu	● synonymous			56433	1	1.77e-5	0	0	0.988
M-12373-A-G		p.Thr13Ala	● missense	Benign		56429	52	9.22e-4	1	1.77e-5	1
M-12375-T-C		p.Thr13Thr	● synonymous			56432	20	3.54e-4	2	3.54e-5	1

Mitochondrial variant page

Single nucleotide variant: M-12451-A-G (GRCh38)

Filter **Pass**

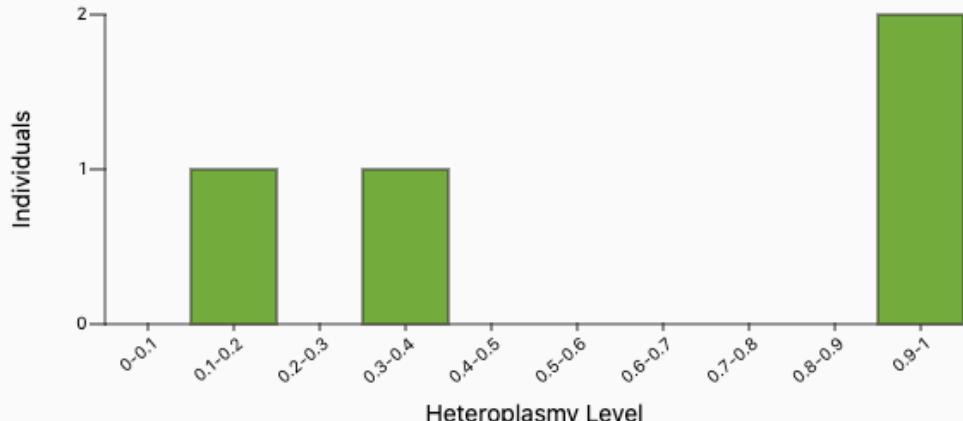
Allele Number: 56429
Homoplasmic AC: 2
Homoplasmic AF: 0.00003544
Heteroplasmic AC: 2
Heteroplasmic AF: 0.00003544
Max Observed Heteroplasmy: 0.9940
Excluded Allele Count: 5
Haplogroup Defining: Yes

Annotations

missense

- **MT-ND5**
- [ENST00000361567.2](#)
HGVS: p.Ile39Val
Polyphen: ● benign
SIFT: ● tolerated_low_confidence

Heteroplasmy Distribution



ClinVar

ClinVar Variation ID: 693446
Conditions: Leigh syndrome
Clinical significance: Likely benign
Review status: criteria provided, single submitter (1 star)
Last evaluated: 2019-10-17

[See submission](#) or find more information on the [ClinVar website](#).

Haplogroup Frequencies ⓘ

Haplogroup	Allele Number	Homoplasmic AC	Homoplasmic AF	Heteroplasmic AC	Heteroplasmic AF
H	14783	0	0.000	1	0.00006765
L2	4724	0	0.000	1	0.0002117
M	1298	1	0.0007704	0	0.000
U	6036	1	0.0001657	0	0.000
Total	26841	2	0.00007451	2	0.00007451

pLoF curations: predictions not the full story

pLoF only Missense only Synonymous only Other only Exomes Genomes SNVs Indels Filter

[Export variants to CSV](#)

Note Only variants located in or within 75 base pairs of a coding exon are shown here. To see variants in UTRs or introns, use the [region view](#).

The table below shows the HGVS consequence and VEP annotation for each variant's most severe consequence across all transcripts in this gene. Cases transcript are denoted with †. To see consequences in a specific transcript, use the [transcript view](#).

Variant ID	Source	HGVS Consequence	VEP Annotation	LoF Curation	Clinical Significance	Flags
1-155206105-C-CT	E	p.Phe386ValfsTer50	● frameshift	Likely LoF		
1-155206261-C-G	E	c.1000-1G>C	● splice acceptor	Likely LoF		
1-155207216-AG-A	E	p.Pro305LeufsTer31	● frameshift	Uncertain	Likely pathogenic	
1-155207370-C-G	E G	c.762-1G>C	● splice acceptor	Likely LoF	Pathogenic	
1-155208051-GA-G	G	p.Ser212HisfsTer19	● frameshift	Likely LoF		
1-155208073-G-A	E	p.Gln205Ter	● stop gained	Likely LoF		
1-155208082-G-A	E G	p.Arg202Ter	● stop gained	Likely LoF		
1-155208089-CAG-C	E	p.Leu199AspfsTer62	● frameshift	Likely LoF	Pathogenic	
1-155208307-C-T	E	c.588+1G>A	● splice donor	Likely LoF		
1-155208442-C-G	E	c.455-1G>C	● splice acceptor	Likely LoF		
1-155209530-CT-C	E	p.Glu111AsnfsTer7	● frameshift	Likely LoF		
1-155209704-GC-G	E	p.Ile95SerfsTer12		This variant was curated as "Likely not LoF". The following factors contributed to this verdict: Homopolymer, Mapping Issue. See variant page for details.		
1-155209741-ACT-A	E	p.Ser81TyrfsTer17				
1-155209764-CA-C	E	p.Gly74ValfsTer17				
1-155209780-C-CG	E	p.Thr69AspfsTer12	● frameshift	Likely not LoF	Pathogenic	
1-155209869-C-T	E	c.116-1G>A	● splice acceptor	Likely LoF		
1-155210420-C-T	E G	c.115+1G>A	● splice donor	Likely LoF	Pathogenic/Likely pathogenic	
1-155210428-C-T	E	p.Trp36Ter	● stop gained	Likely LoF	Pathogenic	pLoF flag
1-155210451-G-GC	E G	p.Leu29AlafsTer18	● frameshift	Likely LoF	Pathogenic	pLoF flag

Manually curated (gnomAD v2)

- All homozygous pLoF variants
 - 28% not LoF
- Heterozygous pLoF variants in some recessive genes
 - 25% not LoF
- Working on pLoF genes in haploinsufficient genes
 - 60% not LoF

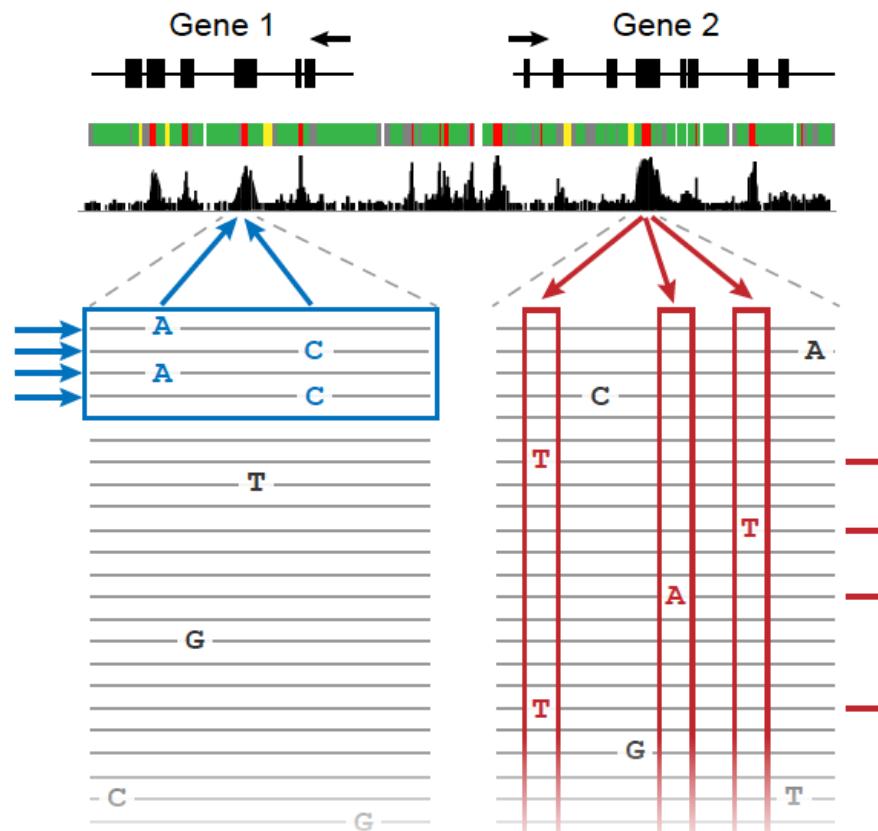
LoF Curation column if the gene has variants curated by the gnomAD team

Genetics as a lever to understand biology



Forward genetics

Find the genetic basis of a phenotype or trait



Discovery of novel disease genes

Reverse genetics

Find phenotypes associated with a particular genetic variant



Phenotyping humans with LoF variants

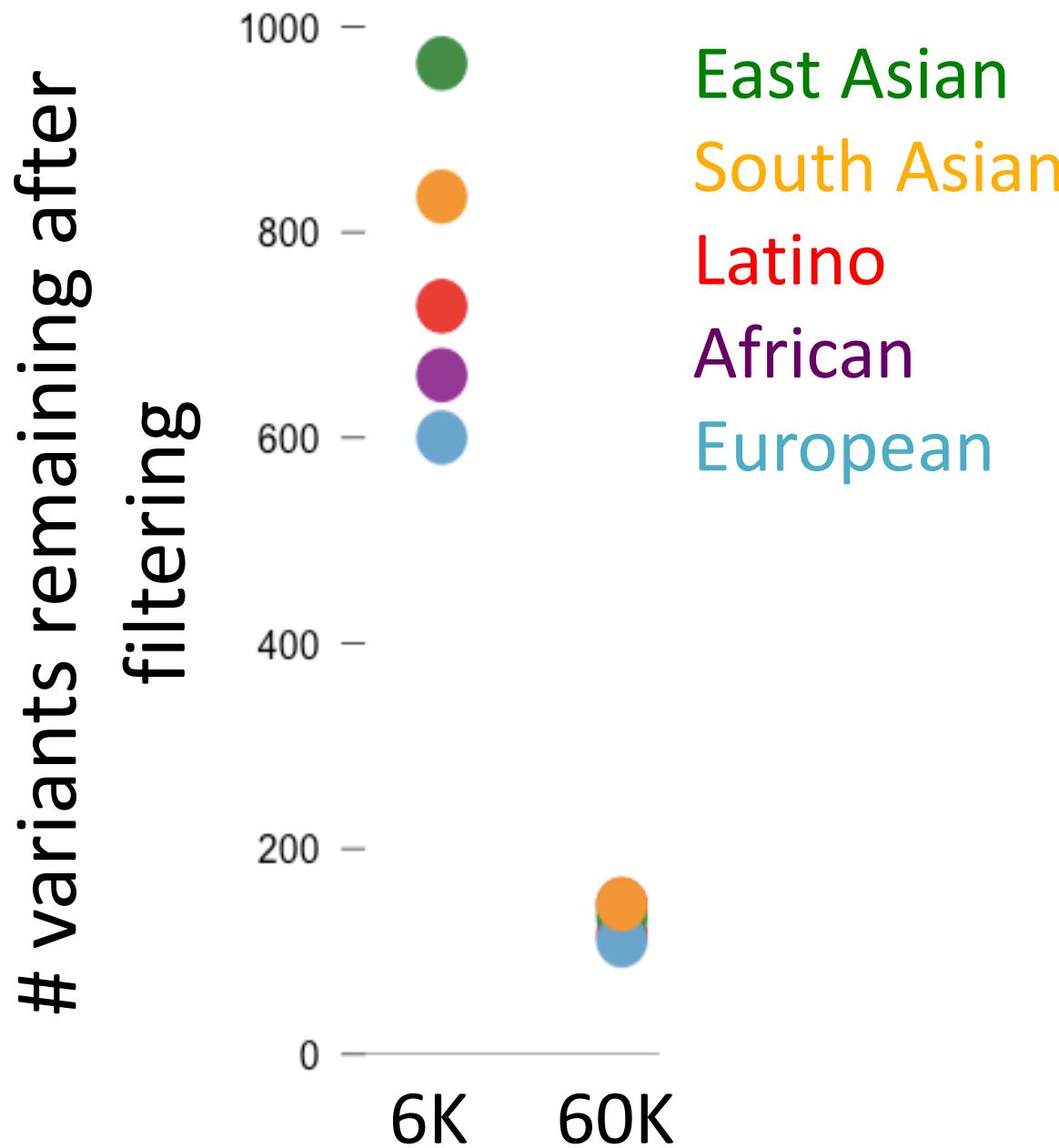
Harnessing the power of allele frequency



VS



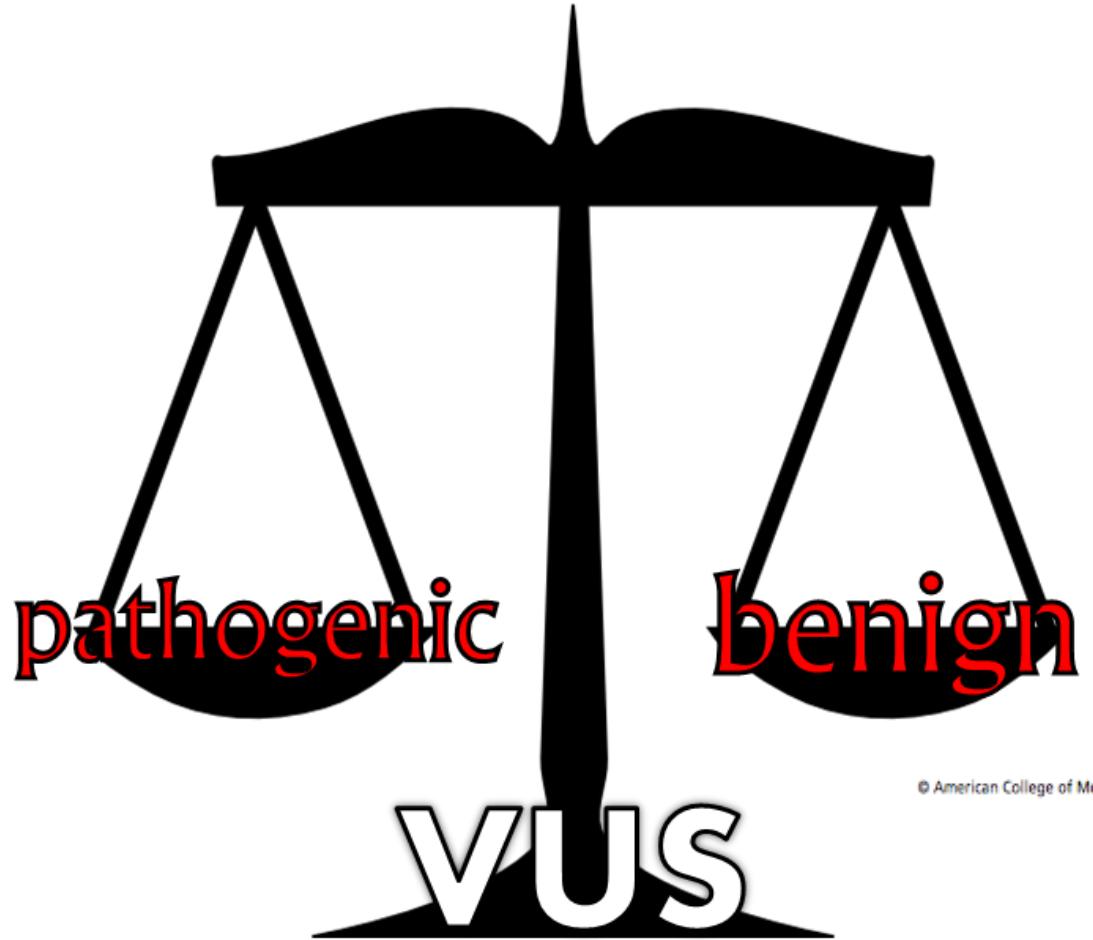
Making sense of one exome requires **tens of thousands** of exomes (or genomes) to reveal rare variants



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



© American College of Medical Genetics and Genomics

Evaluating rare variant pathogenicity

ACMG STANDARDS AND GUIDELINES

Genetics
in Medicine
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

		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			
De novo data				De novo (without paternity & maternity confirmed) PM6	De novo (paternity and maternity confirmed) PS2		
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 <i>trans</i> with a pathogenic variant PM3			

Conclusions

- **Reference population databases** are public resources that are critically important to **evaluate variant rarity**, which is necessary but not sufficient for pathogenicity for rare disease
- These datasets provide multiple tools (constraint, pext) to help prioritize variants to understand human biology
- The power of reference population datasets will increase as they grow in size and diversity
 - gnomAD v4 with >300,000 exomes on GRCh38 anticipated in 2021 or 2022

Downloads

[gnomAD v2](#)[gnomAD v2 liftover](#)[gnomAD v3](#)[ExAC](#)

Available on

- Google Cloud Public Datasets
- Registry of Open Data on AWS
- Azure Open Datasets

Summary

- Variants
- Coverage
- Constraint
- Multi-nucleotide variants (MNVs)
- Proportion expressed across transcripts (pext)
- Structural variants
- Loss-of-function curation results
- Linkage disequilibrium
- Resources

Variants

Note Find out what changed in the latest release in the [gnomAD v2.1.1 README](#).

The variant dataset files below contain all subsets (non-neuro, non-cancer, controls-only, and non-TOPMed).

Exomes

- Sites Hail Table
Show URL for [Google](#) / [Amazon](#) / [Microsoft](#)
Copy URL for [Google](#) / [Amazon](#) / [Microsoft](#)
- All chromosomes sites VCF
58.81 GiB, MD5: f034173bf6e57fbb5e8ce680e95134f2
Download from [Google](#) / [Amazon](#) / [Microsoft](#)
Download TBI from [Google](#) / [Amazon](#) / [Microsoft](#)
- chr1 sites VCF
5.77 GiB, MD5: 9817acdf1d9600efb3355e4cb4b7ee1f
Download from [Google](#) / [Amazon](#) / [Microsoft](#)
Download TBI from [Google](#) / [Amazon](#) / [Microsoft](#)

Genomes

- Sites Hail Table
Show URL for [Google](#) / [Amazon](#) / [Microsoft](#)
Copy URL for [Google](#) / [Amazon](#) / [Microsoft](#)
- All chromosomes sites VCF
460.93 GiB, MD5: e6eadf5ac7b2821b40f350da6e1279a2
Download from [Google](#) / [Amazon](#) / [Microsoft](#)
Download TBI from [Google](#) / [Amazon](#) / [Microsoft](#)
- Exome calling intervals VCF
9.7 GiB, MD5: e5bd69a0f89468149bc3afca78cd5acc
Download from [Google](#) / [Amazon](#) / [Microsoft](#)
Download TBI from [Google](#) / [Amazon](#) / [Microsoft](#)

- Exome and genome data is in separate files
- Exome data is 125K exomes
- Genome data is 15K genomes

gnomAD v3.1 Mitochondrial DNA Variants

November 17, 2020 in [Announcements / Releases](#)

Kristen Laricchia, Sarah E. Calvo

Overview

Mitochondrial DNA (mtDNA) variants for gnomAD are now [available](#) for the first time! We have called mtDNA variants for 56,434 whole genome samples in the v3.1 release. This initial release includes population frequencies for 10,850 unique mtDNA variants defined at more than half of all mtDNA bases. The vast majority of variant calls (98%) are homoplasmic or near homoplasmic, whereas 2% are heteroplasmic. Variation in mitochondrial genomes contributes to many human diseases and has had unique value in the study of human evolutionary genetics. We hope that the addition of mtDNA to gnomAD will enable researchers to better understand the role of mtDNA variation in both health and disease states.

Previous gnomAD callsets have not included mtDNA variants because their properties do not fit the assumptions that we use with our nuclear variant calling pipeline. These properties include:

[Continue reading](#)

gnomAD v3.1

October 29, 2020 in [Announcements / Releases](#)

gnomAD Production Team

Latest posts

[gnomAD v3.1 Mitochondrial DNA Variants](#)
November 17, 2020

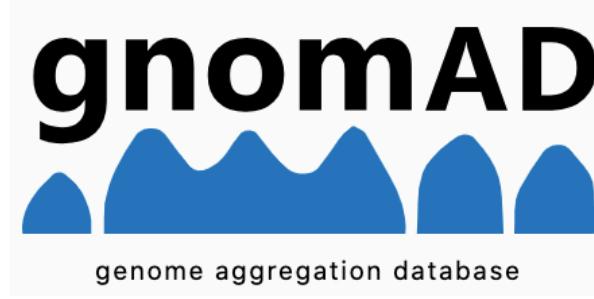
[gnomAD v3.1](#)
October 29, 2020

[gnomAD v3.1 New Content, Methods, Annotations, and Data Availability](#)
October 29, 2020

[Loss-of-Function Curations in gnomAD](#)
October 29, 2020

[Open access to gnomAD data on multiple cloud providers](#)
October 29, 2020

Acknowledgements



Steering Committee

Heidi Rehm (co-PI)
Mark Daly (co-PI)
Daniel MacArthur
Ben Neale
Mike Talkowski
Anne O'Donnell-Luria
Konrad Karczewski
Grace Tiao
Matt Solomonson
Samantha Baxter

Analysis team

Konrad Karczewski
Laurent Francioli
Grace Tiao
Kristen Laricchia
Beryl Cummings
Eric Minikel
Irina Armean
James Ware
Kaitlin Samocha
Nicola Whiffin
Qingbo Wang
Ryan Collins

Ethics team

Andrea Saltzman
Molly Schleicher
Namrata Gupta
Stacey Donnelly

Production team

Grace Tiao
Julia Goodrich
Katherine Chao
Michael Wilson
William Phu
Eric Banks
Charlotte Tolonen
Christopher Llanwarne

David Roazen
Diane Kaplan
Gordon Wade
Jeff Gentry
Jose Soto
Kathleen Tibbetts
Kristian Cibulskis
Laura Gauthier
Louis Bergelson
Miguel Covarrubias
Nikelle Petrillo
Ruchi Munshi
Sam Novod
Thibault Jeandet
Valentin Ruano-Rubio
Yossi Farjoun

Structural Variation

Ryan Collins
Harrison Brand
Konrad Karczewski
Laurent Francioli
Nick Watts
Matthew Solomonson
Xuefang Zhao

Laura Gauthier
Harold Wang
Chelsea Lowther
Mark Walker
Christopher Whelan
Ted Brookings
Ted Sharpe

Jack Fu
Eric Banks
Michael Talkowski



Ben Neale
Cotton Seed
Tim Poterba
Arcturus Wang
Chris Vittal

Browser team

Matthew Solomonson
Nick Watts
Ben Weisburd

Broad Genomics Platform

Stacey Gabriel
Kristen Connolly
Steven Ferriera

<https://gnomad.broadinstitute.org/about>

Contact for ?s: gnomad@broadinstitute.org or odonnell@broadinstitute.org