Using large-scale human genetic variation to inform variant prioritization in neuropsychiatric disorders

Kaitlin E. Samocha

Hurles lab, Wellcome Trust Sanger Institute

ACGS Summer Scientific Meeting 27 June 2017



Critical to determine the subset of variants contributing to disease risk



Well known excess of *de novo* proteintruncating variants (OR \sim 2) in cases with neurodevelopmental disorders.

Modest, but significant, enrichment (~1.4) of *de novo* missense variants.

5,620 cases

5,264 developmental delay / intellectual disability 356 epileptic encephalopathy **2,078 controls**

De Rubeis et al 2014 Iossifov et al 2014 Deciphering Developmental Disorders 2017 EuroEPINOMICS-RES Consortium, EPGP, and Epi4K Consortium 2014 de Ligt et al 2012 Rauch et al 2012 Lelieveld et al 2016

Increasingly large collections of exome sequencing data of reference populations

140,000 -Other 130,000 _atino 120,000 -African Ashkenazi Jewish 110,000 -European 100,000 -South Asian East Asian 90,000 -Individuals 80,000 -70,000 -60,706 60.000 -50.000 -40,000 -30,000 -20,000 -6,503 10,000 -2,504 0 1000ESP ExAC Genomes

Exome Aggregation Consortium (ExAC) 60,706 reference exomes:

- Jointly called and processed
- Unrelated individuals
- Free of individuals with known severe pediatric disease

Increasingly large collections of exome sequencing data of reference populations

141,352 140,000 -Other 130,000 Latino 120,000 -African Ashkenazi Jewish 110,000 -European 100,000 -South Asian East Asian 90,000 -Individuals 80,000 -70.000 -60,706 60,000 -50,000 -40,000 -30,000 -20,000 -6,503 10,000 -2,504 0 1000ESP ExAC gnomAD Genomes

Released in October 2016! gnomAD: 15,136 genomes 126,216 exomes

Lek et al 2016

Challenge of medical genetics: Prioritizing potentially pathogenic variants



- Protein truncating
- Polyphen2 predicted damaging missense
- Presence in reference database



CONSTRAINED × Individual 2 Individual 1 Individual 3 Individual 5 Individual 4 Individual 6 Т Ι Μ Е \star

How few mutations indicate constraint?



We need a way to determine if the number of observed variants is significantly different from expectation



Mutational model accurately predicts synonymous variation

- We used our mutational model to predict the expected number of variants in the ~61k individuals in the ExAC dataset
- Extracted rare (MAF < 0.1%) variants as a comparison and found a high correlation for synonymous ($r^2 = 0.96$)



The ~61,000 exomes allowed us to:

- 1. Evaluate constraint against protein-truncating variants
- 2. Investigate the missense constraint of regions within genes

pLI: the probability of being loss-of-function intolerant

- Genes that are extremely loss-of-function (LoF) intolerant will be depleted of such variation in a reference population
- The proposed mechanism of this intolerance is haploinsufficiency
- Created an EM-based metric that broadly divides genes into two categories ("likely LoF intolerant" and "not likely LoF intolerant")

The creation of pLI

We propose a model based on Mendelian modes of inheritance:

Tolerant of loss of both copies Tolerant of loss of a single copy Intolerant of loss of a single copy

The creation of pLI

We propose a model based on Mendelian modes of inheritance:



Gene: GRIN2B



The probability of being loss-of-function intolerant (pLI) shows the expected contrast between gene lists

Severe haploinsufficient (n = 41)Moderate haploinsufficient (n = 72)Mild haploinsufficient (n = 58)Essential in culture (n = 280)Dominant disease (n = 693)All genes Н (n = 18,225)**Recessive disease** (n = 1, 167)Olfactory (n = 351)0.1 0.2 0.4 0.50.70.3 0.6 0.8 0 Anne O'Donnell-Luria Fraction of gene set that is highly LoF-intolerant Emma Pierce-Hoffman $(pLI \ge 0.9)$ ClinGen; OMIM

Lek et al 2016

1

0.9

The probability of being loss-of-function intolerant (pLI) shows the expected contrast between gene lists



The probability of being loss-of-function intolerant (pLI) shows the expected contrast between gene lists



What is pLI truly capturing?

Genes that, when disrupted, cause conditions that are:





Anne O'Donnell-Luria Emma Pierce-Hoffman ClinGen; OMIM

What is pLI truly capturing?

Genes that, when disrupted, cause conditions that are:

≻Severe

Dominant (or haploinsufficient)



What is pLI truly capturing?

Genes that, when disrupted, cause conditions that are:

≻Severe

>Dominant (or haploinsufficient)

>Early onset

BRCA2 has 52% of expected protein truncating variants, giving it a pLI ~ 0

Applying pLI to *de novo* proteintruncating variants



5,620 cases

5,264 developmental delay/ intellectual disability 356 epileptic encephalopathy **2,078 controls** Combining variant and gene level evidence to identify a high impact subset of *de novo* PTVs





Unaffected

Kosmicki et al 2017

Combining variant and gene level evidence to identify a high impact subset of *de novo* PTVs



The ~61,000 exomes allowed us to:

- 1. Evaluate constraint against protein-truncating variants
- 2. Investigate the missense constraint of regions within genes

We expect that for some genes, only sections of them will truly be missense constrained

Some genes have regions of missense constraint

Example gene: 401 missense variants expected and 199 (~50%) observed.



Some genes have regions of missense constraint

Example gene: 401 missense variants expected and 199 (~50%) observed.





base pairs







MPC: combining local constraint with variant level information

Missense badness

amino acid substitution deleteriousness metric similar to BLOSUM/Grantham

PolyPhen-2

missense pathogenicity metric

Constraint

missense depletion of region or gene

Missense variants in ExAC with MAF > 0.1% (n = 82,932)



Does MPC help differentiate likely benign from likely pathogenic *de novo* variants?



5,620 cases

5,264 developmental delay/ intellectual disability 356 epileptic encephalopathy 2,078 controls









Comparing MPC to other metrics

Jointly ranked *de novo* missense variants from cases and controls by each metric and evaluated the fraction of case variants in the top 10%. The total proportion of case variants is 0.8 (5113/6382), so a metric with no predictive value would match this overall rate.

	MPC	M-CAP	CADD	PolyPhen-2
Fraction of top 10% from cases	0.95	0.93	0.86	0.85
Odds ratio	5.43	3.52	1.58	1.44
p-value	1.48×10^{-28}	4.35x10 ⁻²⁰	1.46x10 ⁻⁴	2.66x10 ⁻³

Jagadeesh et al 2016; Kircher et al 2014; Adzhubei et al 2010

Availability of the regional constraint data and method

Preprint describing the regional constraint work is on bioRxiv:



doi: https://doi.org/10.1101/148353

And you can download MPC scores for all possible missense variants in ~18k canonical transcripts from the ExAC FTP:

ftp://ftp.broadinstitute.org/pub/ExAC_release/release1/regional_missense_con straint/ Created a mutational model to predict the expected amount of variation in a reference population:

- Identified ~3k genes that are extremely intolerant of lossof-function variation
- The most missense depleted regions of are enriched for pathogenic variants and *de novo* missense variants in cases with a neurodevelopmental disorder



Acknowledgements

Mark Daly Benjamin Neale Daniel MacArthur Jack Kosmicki Konrad Karczewski Monkol Lek Anne O'Donnell-Luria Elise Robinson Eric Minikel Emma Pierce-Hoffman

Exome Aggregation Consortium (ExAC) <u>http://exac.broadinstitute.org/about</u>

genome Aggregation Database (gnomAD) http://gnomad.broadinstitute.org/about

XXX

