

Algorithms to elucidate causative gene mutations in rare disorders

PHEVOR & Phen-Gen

Journal Club

Marie-Angela Wulf

9.9.2014

Overview

- 1 Motivation
- 2 Background
 - short repetition in genetics
 - how to find the causative gene?
 - algorithms that are currently available
- 3 PHEVOR
 - Methods
 - Results
- 4 Phen-Gen
 - Phenotypic Prediction
 - Genotypic Prediction
 - Results
 - Comparison of the two
- 5 Discussion & Conclusion

- is there a gene mutation leading to disease?
- if so, which one?
- is it inherited?
- is there a risk of replication in future siblings?
- possibility of pre-natal diagnostics?

1

¹ <http://saalfeld.otz.de/web/lokal/leben/detail/-/specific/Familie-aus-Saalfeld-mit-behinderter-Tochter>

Motivation



- is there a gene mutation leading to disease?
- if so, which one?
- is it inherited?
- is there a risk of replication in future siblings?
- possibility of pre-natal diagnostics?
- Problem: very rare diseases
=> very low number of cases
- only in approx. 50% of cases, the causal gene is known
- gap between sequencing and analysis capabilities

2

²[http://saalfeld.otz.de/web/lokal/leben/detail/-/specific/Familie-aus-Saalfeld-mit-behinderter-Tochter-](http://saalfeld.otz.de/web/lokal/leben/detail/-/specific/Familie-aus-Saalfeld-mit-behinderter-Tochter-braucht-Hilfe-291757532)

short repetition in genetics

Exome coding region of the genome

start-loss mutation in the ATG start codon that prevents the start of translation

stop-gain mutation that induces a new stop-codon

stop-loss loss of the normal stop codon, leading to larger transcript

splice site where introns are cut out and exons are put together for the final mRNA

nonsynonymous mutation leads to incorporation of a different amino acid

synonymous mutation leads to incorporation of the same amino acid

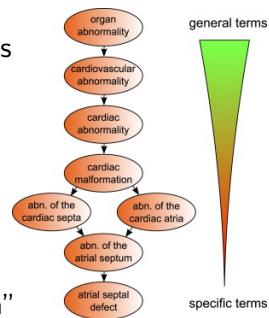
indel small insertion or deletion

Human Phenotype Ontology (HPO)

Ontology computational representation of a domain of knowledge based upon a controlled, standardized vocabulary for describing entities and the semantic relationships between them

Annotation is the process of assigning ontology terms (concepts) for the description of objects

- terms are related to parent terms by “is a” relationship
- different parent terms possible



Human Phenotype Ontology (HPO)

- introduced by Robinson et al in 2008
- standardized vocabulary of phenotypic abnormalities encountered in human disease
- Each term in the HPO describes a phenotypic abnormality
- ontology terms are used for annotation of diseases
- almost 50,000 annotations to 4813 diseases are provided

how to find the causative gene?

- sequencing
- comparison to reference data set
- but: also healthy individuals have about 100 genuine loss-of-function variants with about 20 genes completely inactivated
- \Rightarrow need of a good post-sequencing processing step in order to find the causative mutation

algorithms that are currently available

- VAAST (2011)
- eXtasy (2013)
- FunSeq (2013)
- PHIVE (2014)
- PHEVOR (2014)
- Phen-Gen (2014)

algorithms that are currently available

- VAAST (2011)
- eXtasy (2013)
- FunSeq (2013)
- PHIVE (2014)
- *PHEVOR (2014)*
- *Phen-Gen (2014)*

principles of variant-prioritization tools at the example of VAAST

- ① Sequencing
- ② Comparison with reference genome or exome sequence => variants
- ③ ranking of the variants according to different criteria
 - ① frequency of the variant in reference genomes
 - ② predicted severity of amino acid changes
 - ③ location of the variant in the genome
 - ④ inheritance and penetrance patterns

This leads to causal gene top rank in 62% of cases for dominant disorders and 78% for recessive disorders.

PHEVOR

Phevor Combines Multiple Biomedical Ontologies for Accurate Identification of Disease-Causing Alleles in Single Individuals and Small Nuclear Families

Marc V. Singleton,¹ Stephen L. Guthery,⁵ Karl V. Voelkerding,^{3,4} Karin Chen,⁶ Brett Kennedy,¹ Rebecca L. Margraf,⁴ Jacob Durtschi,⁴ Karen Eilbeck,^{2,7} Martin G. Reese,⁸ Lynn B. Jorde,^{1,2} Chad D. Huff,⁹ and Mark Yandell^{1,2,*}

¹Department of Human Genetics, Eccles Institute of Human Genetics, University of Utah School of Medicine, Salt Lake City, UT 84112, USA; ²Utah Science, Technology, and Research Center for Genetic Discovery, University of Utah, Salt Lake City, UT 84112, USA; ³Department of Pathology, University of Utah, Salt Lake City, UT 84112, USA; ⁴ARUP Institute for Clinical and Experimental Pathology, 500 Chipeta Way, Salt Lake City, UT 84108, USA; ⁵Division of Pediatric Gastroenterology, Hepatology, and Nutrition, Department of Pediatrics, University of Utah, Salt Lake City, UT 84112, USA; ⁶Division of Allergy, Immunology, and Rheumatology, Department of Pediatrics, University of Utah, Salt Lake City, UT 84112, USA; ⁷Department of Biomedical Informatics, University of Utah, Salt Lake City, UT 84112, USA; ⁸Omicron Inc., 1625 Clay Street, Oakland, CA 94612, USA; ⁹Department of Epidemiology, The University of Texas MD Anderson Cancer Center, P.O. Box 301439, Houston, TX 77230, USA

The American Journal of Human Genetics 94, 599–610, April 3, 2014

Introduction

- PHEVOR = Phenotype Driven Variant Ontological Re-ranking tool
- Idea: combine the outputs of variant-prioritization tools with knowledge resident in biomedical ontologies ^a

Phenotype Description

- Phenomizer Web Tool
- Human Phenotype Ontology Terms
- Gene Ontology Terms
- Disease Ontology Terms

Phevor

Re-Ranked Candidate Gene List

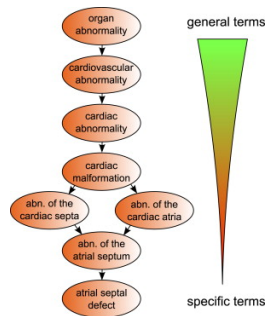
- SIFT
- AnnoVar
- PhastCons
- VAAST

Variant Prioritization

^aRobinson et al, 2008, Smith et al, 2012, Schriml et al, 2012, Ashburner et al, 2000

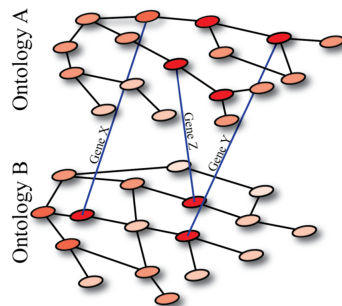
Workflow

- 1 user feeds the software with a list of phenotypic terms from an ontology (or a gene list)
- 2 Phevor saves the annotated genes in an internal list. If there are no annotated genes, Phevor iterates towards the root of the ontology and saves the annotated genes of the next nodes
- 3 the gene list is then fed into the other ontologies

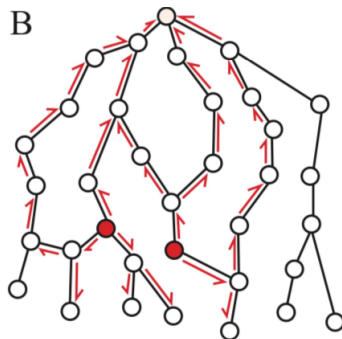
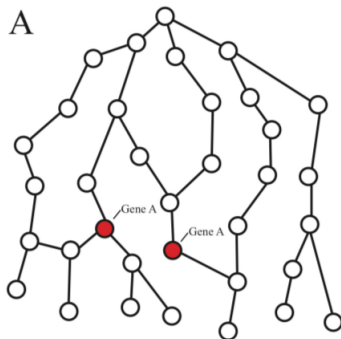


Workflow

- 1 user feeds the software with a list of phenotypic terms from an ontology (or a gene list)
- 2 Phevor saves the annotated genes in an internal list. If there are no annotated genes, Phevor iterates towards the root of the ontology and saves the annotated genes of the next nodes
- 3 the gene list is then fed into the other ontologies

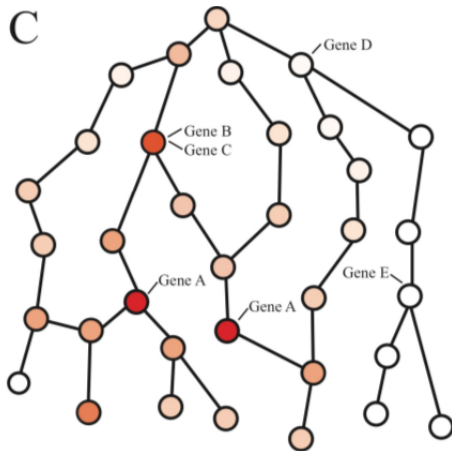


Ontological Propagation



each time an edge is crossed to a neighboring node, the current value of the previous node is divided by 2.

Ontological Propagation



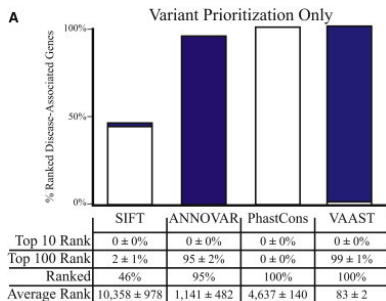
- Phevor renormalizes each node's value to between 0 and 1 by dividing it by the sum of all node scores in the ontology.
- $\text{Score}(\text{Gene})$ in a given ontology is the maximum score it achieved in any node
- final score is the sum of all the scores the gene achieved in the different ontologies
- normalized again to a value between 0 and 1

Combining Ontologies and Variant Data

- disease association score for each gene: $D_G = (1 - V_G) \times N_G$
where N_G is the percentile rank of the renormalized gene sum score as derived from the ontological combination and propagation procedures and V_G is the gene's percentile rank provided by the external variant-prioritization or search tool
- H_G , a second score summarizing the weight of evidence that the gene is not involved with the individual's illness:
$$H_G = V_G \times (1 - N_G)$$
- Phevor score: $S_G = \log \frac{D_G}{H_G}$

Results: gene known

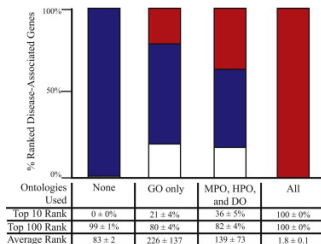
- two copies of a known disease-causing allele were inserted into an exome
- the process was repeated 100 times for 100 different genes with known disease associations



- two copies of a known disease-causing allele were inserted into an exome
- the process was repeated 100 times for 100 different genes with known disease associations



Results: gene not known



only VAAST genetic input

none no annotated genes

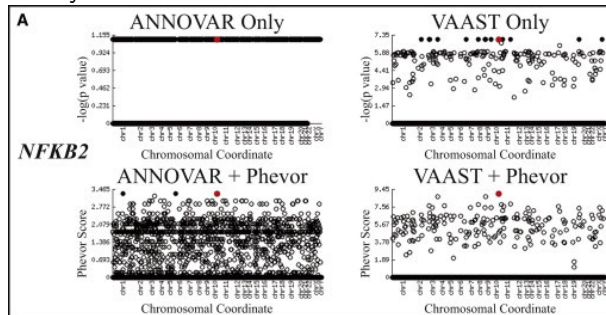
GO only annotated genes
only in GO

MPO, HPO & DO annotated
genes removed only
in GO

all annotated genes
there

Results: Application in real patients 1

family affected by autosomal-dominant, early-onset hypogammaglobulinemia with variable autoimmune features and adrenal insufficiency



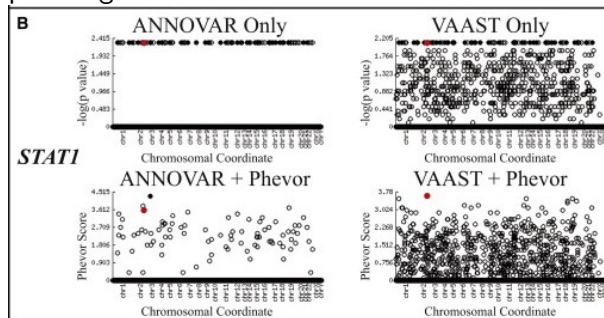
3

Results: Application in real patients 2

12-year-old male with severe diarrhea in the context of intestinal inflammation, total villous atrophy, and hypothyroidism

=> X-linked immunodysregulation, polyendocrinopathy, and enteropathy (IPEX syndrome)

BUT clinical sequencing of genes associated with IPEX revealed no pathologic variants



4

Summary

- PHEVOR ameliorates the performance of sequencing-based ranking algorithms
- can identify dominant mutations in a gene not previously associated with the disease phenotype in the ontologies (NFKB2) or a de novo dominant allele located in an known disease-associated gene (STAT1) and producing an atypical phenotype.
- if the gene is not annotated in the ontologies, the performance drops

Phen-Gen

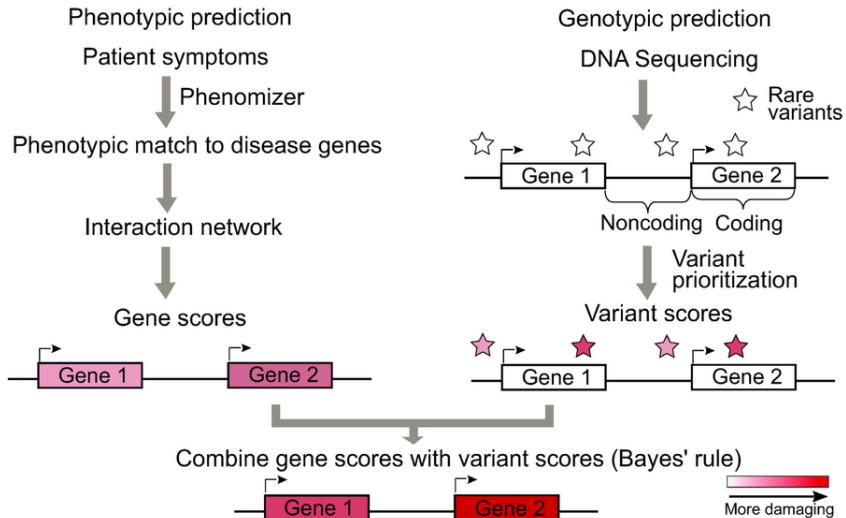
Phen-Gen: combining phenotype and genotype to analyze rare disorders

Asif Javed, Saloni Agrawal & Pauline C Ng

Computational and Systems Biology Group, Genome Institute of Singapore, Agency for Science, Technology and Research, Singapore.

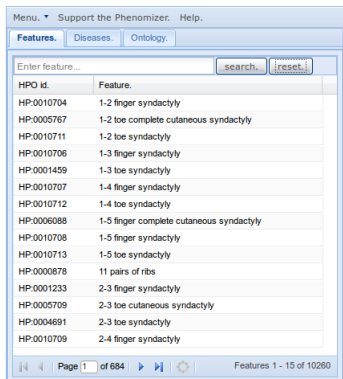
NATURE METHODS | VOL.11 NO.9 | SEPTEMBER 2014 | 935

overall workflow



Phenomizer

- introduced by Kohler et al, 2009
- Aim: match the patient's symptomatology to the list of known disorders and estimate the significance of each disease match



Phenomizer

- introduced by Kohler et al, 2009
- Aim: match the patient's symptomatology to the list of known disorders and estimate the significance of each disease match

Menu. ▾ Support the Phenomizer. Help. The Phenomizer

Features. Diseases. Ontology.

Enter feature...

HPO id.	Feature.
HP:0010704	1-2 finger syndactyly
HP:0005767	1-2 toe complete cutaneous syndactyly
HP:0010711	1-2 toe syndactyly
HP:0010706	1-3 finger syndactyly
HP:0001459	1-3 toe syndactyly
HP:0010707	1-4 finger syndactyly
HP:0010712	1-4 toe syndactyly
HP:0006088	1-5 finger complete cutaneous syndactyly
HP:0010708	1-5 finger syndactyly
HP:0010713	1-5 toe syndactyly
HP:0000878	11 pairs of ribs
HP:0001233	2-3 finger syndactyly
HP:0005709	2-3 toe cutaneous syndactyly
HP:0004691	2-3 toe syndactyly
HP:0010709	2-4 finger syndactyly

Page 1 of 684 Features 1 - 15 of 10260

Patient's Features. **Diagnosis.**

Algorithm: resnik (Unsymmetric). 2 Features.

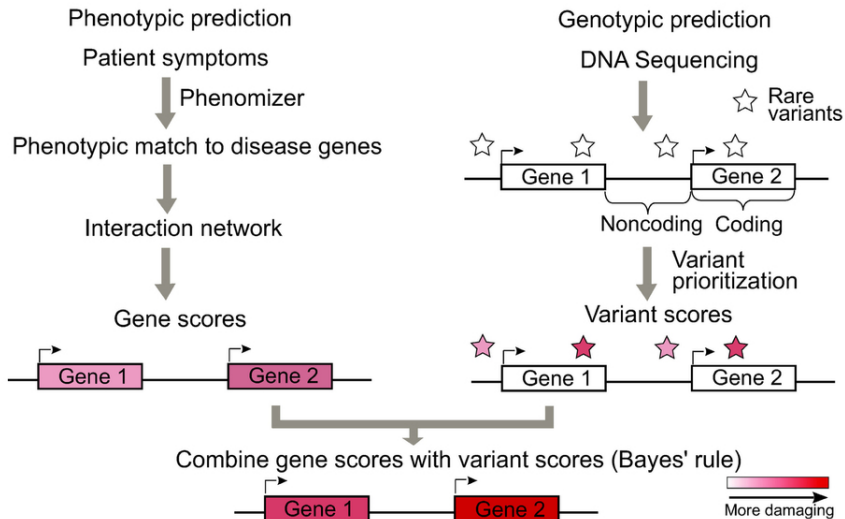
<input type="checkbox"/>	p-value. ▲	Disease Id.	Disease name.	Genes.
<input checked="" type="checkbox"/>	0.0736	OMIM:12...	CREUTZFELDT-JAKOB DISEASE	PRNP
<input checked="" type="checkbox"/>	0.0736	OMIM:60...	#603516 SPINOCEREBELLAR ATAXIA 10; SCA...	ATXN10
<input checked="" type="checkbox"/>	0.0736	OMIM:60...	#607136 SPINOCEREBELLAR ATAXIA 17; SCA...	TBP
<input checked="" type="checkbox"/>	0.0736	OMIM:13...	#137440 GERSTMANN-STRAUSSLER DISEA...	PRNP
<input checked="" type="checkbox"/>	0.0736	OMIM:30...	#300623 FRAGILE X TREMOR/ATAXIA SYND...	FMR1
<input checked="" type="checkbox"/>	0.0736	OMIM:24...	#249900 METACHROMATIC LEUKODYSTRO...	PSAP
<input checked="" type="checkbox"/>	0.0736	OMIM:20...	#208920 ATAXIA, EARLY-ONSET, WITH OCUL...	APTX
<input checked="" type="checkbox"/>	0.0736	OMIM:60...	SPINOCEREBELLAR ATAXIA 14	PRKCG
<input checked="" type="checkbox"/>	0.0736	OMIM:61...	#610217 NEURODEGENERATION WITH BRAI...	PLA2G6
<input checked="" type="checkbox"/>	0.0736	OMIM:61...	#612075 MITOCHONDRIAL DNA DEPLETION ...	RRM2B
<input checked="" type="checkbox"/>	0.0736	OMIM:60...	#604121 CEREBELLAR ATAXIA, DEAFNESS, A...	DNMT1
<input checked="" type="checkbox"/>	0.0736	OMIM:30...	301840 SPINOCEREBELLAR ATAXIA, X-LINK...	
<input checked="" type="checkbox"/>	0.0736	OMIM:11...	DEMENTIA, FAMILIAL DANISH	ITM2B
<input checked="" type="checkbox"/>	0.0736	OMIM:16...	164700 OLIVOPONTOCEREBELLAR ATROPH...	
<input checked="" type="checkbox"/>	0.0736	OMIM:61...	#615362 CEROID LIPOFUSCINOSIS, NEURO...	CTSF

Page 1 of 249 Improve Differential Diagnosis. Download Result

Phenomizer in Phen-Gen

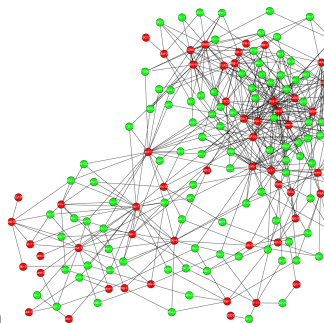
- translation of p-value to disease-probability
- Each disease probability is assigned to all genes implicated in the database for that disorder. If no genes are currently known, the probability is distributed uniformly across all genes.

overall workflow



Interaction Network

- Interaction: ability of two genes to affect the same underlying biology and thus lead to the same (or similar) disorders and symptoms
- Interaction data from REACTOME, KEGG, NCI-Nature, BioGRID, STRING, Gene Ontology domain annotations, COXPRESdb
- integrated into a gene-gene interaction network
- edge weights proportional to confidence in interaction
- 920,898 interactions in total

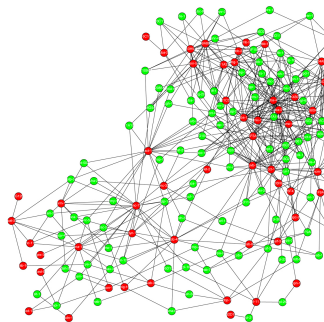


5

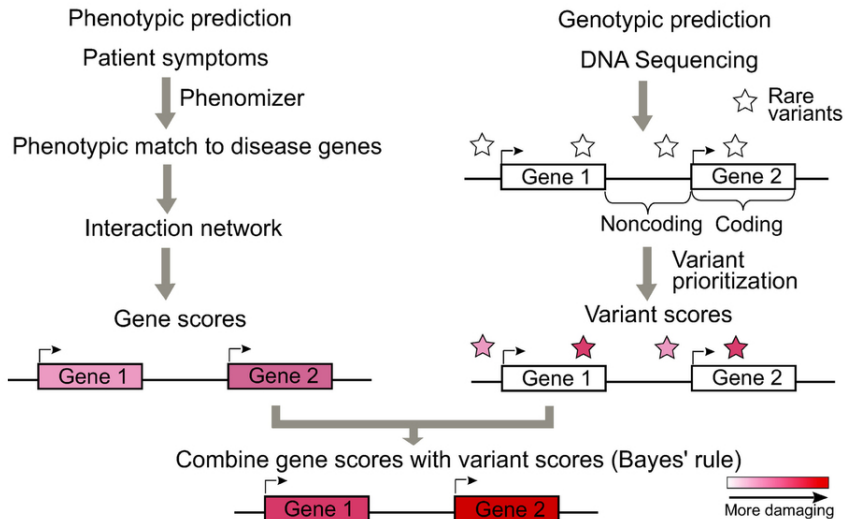
⁵ Matthews et al, 2009, Kanehisa et al, 2012, Schaefer et al, 2009, Stark et al, 2006, Franceschini et al, 2013,

Interaction Network

- random walk with restart on the gene-gene interaction network.
 - Each gene selected as starting point of the walk with probability P_{Gene} with restart probability of 90%
 - \Rightarrow 90% of the 'guilt', retained by the initial gene & 10% of it permeated across its known associates (with stronger evidence of interaction leading to higher probabilities)



overall workflow



Gene List

- four popular gene lists (Consensus CDS, RefSeq, Ensembl and UCSC known)
- Incomplete transcripts discarded
- only protein-coding genes considered
- Alternatively spliced transcripts reported in the same data source were assigned to the same gene
- Additionally, two genes were merged if any of their isoforms exhibited greater than 95% identity in the exon sequence in the same orientation
- The final gene list comprised 26,803 protein-coding genes
- key question: is the reported mutation harmful or not?

Variant Prediction - exome-centric approach

- each variant in the patient's genome (or exome) evaluated if
 - within the coding regions in a reported transcript of the gene list
 - within the splice-site definition of the intron-exon boundary
- coding variants subcategorized as start-loss, stop-gain, stop-loss, splice site, nonsynonymous, synonymous or indel.
- Each variant is assigned a probability of deleteriousness based on its highest estimated damaging impact.

Variant Prediction - exome-centric approach

Nonsynonymous variants For prediction of damaging impact, two commonly used algorithms were combined: SIFT and PolyPhen-2 => estimate the probability of a nonsynonymous mutation

Variant Prediction - exome-centric approach

Nonsynonymous variants For prediction of damaging impact, two commonly used algorithms were combined: SIFT and PolyPhen-2 => estimate the probability of a nonsynonymous mutation

Splice-site variants Mutations in the 8-bp locus surrounding the donor site or the 3-bp locus surrounding the acceptor site were considered splice disrupting

Variant Prediction

Start-loss, stop-gain & stop-loss variants are highly deleterious and assumed to be damaging

Variant Prediction

Start-loss, stop-gain & stop-loss variants are highly deleterious and assumed to be damaging

genic indels the prior probability of a genic indel being damaging was estimated to be 0.0787 based on analysis of a reference data set.

frameshift indels are assigned the probability of deleteriousness on the basis of the empirical distribution of its positive and neutral set combined nonframeshift indel predictions are further refined by incorporating the importance of the affected locus (based on its tolerance to single-nucleotide mutations)

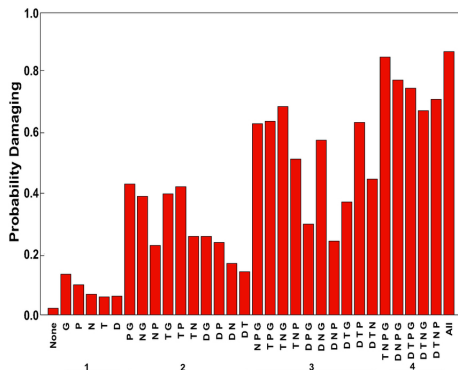
Variant Prediction - genome-centric approach

Genomic variants

aim: to estimate the putative functional role of all genomic variants, esp. those lying outside of the exome

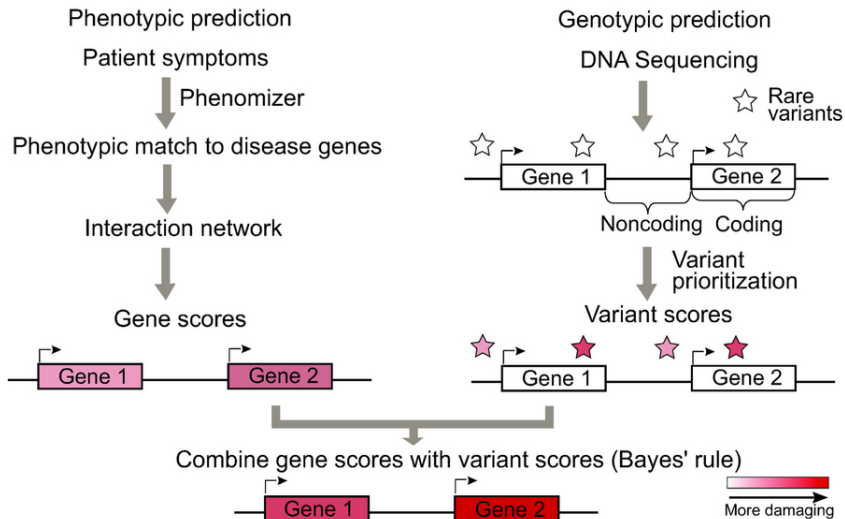
- conservation (using GERP++ and PhyloP)
- putative regulatory interactions (using ENCODE)
- proximity to genes (including coding sequence, UTRs and 70 bp at the start and end of each intron)
- => annotate the variants in order to estimate a locus's susceptibility for disease

Variant Prediction- genome-centric approach



G: GERP++, P: PhyloP, N: near-genic, T: transcription factor binding sites, and D: DNase hypersensitive sites

overall workflow

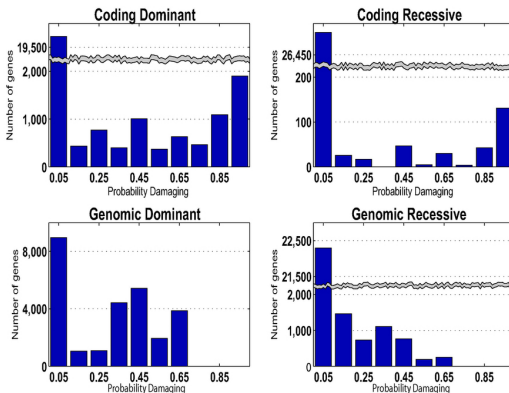


Pooling

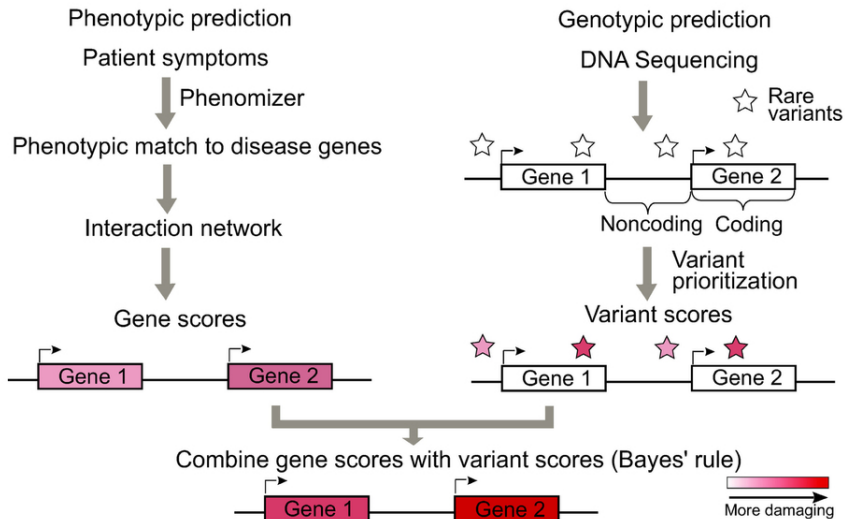
- The estimated loss of function at the genic level is predicted by pooling damaging variants within each gene considering the most damaging predicted variants for maternally and paternally inherited chromosomes.
- The inheritance pattern of the variants is compared against the disease inheritance pattern.
- Common variants are omitted from this pooling.
- only genes harboring variants that exceed the 99th percentile of the corresponding null distribution are considered for downstream analysis.

Pooling - null distribution

null distribution of predictor for each gene using the samples from 1000 healthy individuals (top 1 percentile of damaging variants)



overall workflow



Combined Prediction

the probability that gene j is disease causal is given by

$$\frac{P_G(j) \times P_P(j)}{P_G(j) \times P_P(j) + (1 - P_G(j)) \times (1 - P_P(j))}$$

where $P_G(j)$ and $P_P(j)$ are the probabilities that gene j is disease causal on the basis of genotype data and phenotype data, respectively

Phen-Gen online

Phen-Gen

[HOME](#)[ABOUT](#)[INSTRUCTIONS](#)[DEMO](#)[RESULTS](#)[DOWNLOAD](#)[CONTACT US](#)

Phen-Gen online

INPUT FILES ^[?] (*required)

Variants:* Keine Datei ausgewählt.

In [VCF](#) [Sample VCF](#)

Phenotypes:* Keine Datei ausgewählt.

In a text file with 1 [HPO](#) term per line. [Sample phenotype file](#)

Pedigree:* Keine Datei ausgewählt.

In [PED](#) format. [Sample PED](#)

RUN PARAMETERS ^[?]

Disease inheritance pattern:

☒ Recessive

☐ Dominant

Type of prediction:

☒ Coding

☐ Genomic

Discard de novo mutations:

☒ Yes

☐ No

Stringency:

Only digits (0-9)

Your e-mail address:

To get notified when the job is done.

Not sure what to do? See [instructions](#) and a [demo](#).

Approximate run time: 15-30 minutes

Results: Evaluation of exome-predictor with in-silico patients

- a patient's genetic data were created by adding a disease-causal mutation into a healthy individual's genome or exome
- phenotypic data was generated on the basis of reported disease symptoms

	Dominant			Recessive		
	Top gene (%)	Top 5 genes (%)	Top 10 genes (%)	Top gene (%)	Top 5 genes (%)	Top 10 genes (%)
Combined (missense, nonsense, splice site and indels)						
Phen-Gen	81	88	91	96	97	97
Genotype only	0	3	13	72	97	97
Phenotype only	31	69	74	24	51	68
Known	92	97	97	96	97	97
Unknown	43	61	73	92	96	96

Results: Evaluation of genomic predictor with in-silico patients

Performance if regulatory elements are affected

	Dominant			Recessive			Combined (dominant + recessive)		
	Top gene (%)	Top 5 genes (%)	Top 10 genes (%)	Top gene (%)	Top 5 genes (%)	Top 10 genes (%)	Top gene (%)	Top 5 genes (%)	Top 10 genes (%)
Regulatory									
Phen-Gen	58	64	65	40	61	72	49	62	69
Genotype only	1	2	2	34	51	62	19	29	34
Phenotype only	31	43	54	5	11	26	17	26	39

Results: Evaluation with real patients

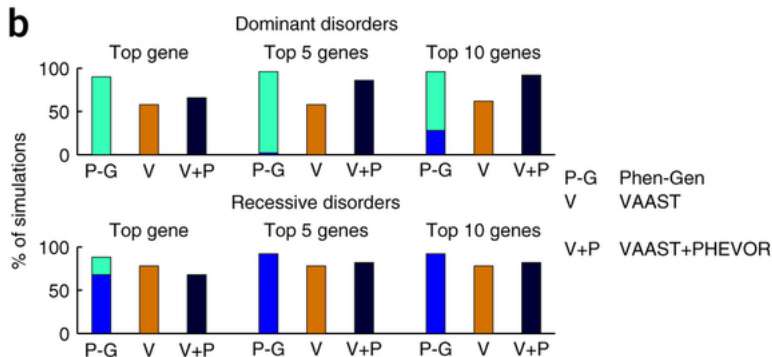
Phen-Gen applied to eleven families with variants implicated in recessive or X-mode inheritance

Trio ID	Gene	Classification	Phen-Gen's Rank	
			Inclusive of dbSNP MAF<1%	Filtration criteria from ref. 5
4	<i>FANCB</i>	Unknown	18	3
4	<i>PDHA1</i>	Known	5	2
4	<i>GUCY2F</i>	Unknown	4	1
16	<i>ENOX2</i>	Unknown	2	1
18	<i>ARHGEF9</i>	Known	8	1
25	<i>GPM6B</i>	Unknown	7	1
41	<i>ARHGEF9</i>	Known	1	1
42	<i>DDX26B</i>	Unknown	13	3
72	<i>PDZD11</i>	Unknown	3	3
93	<i>TRPC5</i>	Candidate	4	1
12	<i>SYCP2L</i>	Unknown	2	1
12	<i>VPS13B</i>	Known	14	4
12	<i>C8orf59</i>	Unknown	4	2
12	<i>PRUNE2</i>	Unknown	7	3
24	<i>PCNT</i>	Known	7	2
70	<i>IQGAP2</i>	Unknown	6	1

Filtration criteria: all

mutations reported in the Single Nucleotide Polymorphism Database (dbSNP) or observed in the in-house data set were removed from further evaluation

Comparison of Phen-Gen with PHEVOR



Phen-Gen outperforms VAAST by 19-21% and PHEVOR by 20-26% when considering the first rank

Summary

- Phen-Gen top-ranks the causal gene in 49% (non-coding variants) to 97% (coding recessive variants)
- it outperforms other existing methods
- fast (15-30min)
- open source license
- clinical use needs to be confirmed
- a better performance with non-coding mutations would be nice

Discussion & Conclusion

Discussion

- need of very careful clinical diagnostics
- clinical value needs to be confirmed
- performance esp. in non-coding mutations might be ameliorated
- rely on further development of other existing tools

Discussion & Conclusion

Discussion

- need of very careful clinical diagnostics
- clinical value needs to be confirmed
- performance esp. in non-coding mutations might be ameliorated
- rely on further development of other existing tools

Conclusion

such algorithms might facilitate

- diagnosis of rare genetic diseases
- genetic counseling
- identification of candidate genes for further research
- drug development

List of Links

HPO: <http://www.human-phenotype-ontology.org/>

PHEVOR: <http://weatherby.genetics.utah.edu/cgi-bin/Phevor/PhevorWeb.html>

VAAST: <http://www.yandell-lab.org/software/vaast.html>

Phenomizer: <http://compbio.charite.de/phenomizer/> or in App-/Play- Store

Phen-Gen: <http://phen-gen.org/>

Questions?



Thank you!



Appendix: Bayes' Rule

- Bayes' Theorem: Rule for conditional probability i.e. the probability of A given B
 - calculation: $P(A | B) = \left(\frac{P(A) \times P(B|A)}{P(B)} \right)$
- Bayes' Rule: relates the odds of event A_1 to the odds of event A_2 , before (prior to) and after (posterior to) conditioning on another event B
 - calculation: $\frac{P(A_1|B)}{P(A_2|B)} = \left(\frac{P(A_1) \times P(B|A_1)}{P(A_2) \times P(B|A_2)} \right)$

Appendix: Phenomizer in detail

- uses HPO
- Calculation of the similarity score:
 - calculation of Information Content:

$$IC = -\log\left(\frac{\text{number of diseases having the term}}{\text{number of all diseases}}\right)$$
 - \Rightarrow the more specific the term, the higher the IC
 - similarity score: $sim = avg [\sum IC(MICA)]$ where MICA=most informative common ancestor
 - weighting function for additional features of the disease not covered by the entered terms
- Calculation of p-value:
 - based on the distribution of similarity scores that is obtained by randomly choosing combinations
 - If a given score is only rarely obtained by chance, then we consider it to be statistically significant.
 - Monte Carlo random sampling and corrected for multiple testing by the method of Benjamini and Hochberg