MutationTaster & RegulationSpotter

Pathogenicity Prediction of Sequence Variants: Past, Present and Future

Dr. rer. nat. Jana Marie Schwarz

Klinik für Pädiatrie m. S. Neurologie Exzellenzcluster NeuroCure Charité Universitätsmedizin Berlin

AG Translational Genomics

07.11.2017

AG Translational Genomics

Focus on rare (mendelian) diseases in children

muscular, neurological and neuro-muscular system

Next Generation Sequencing (NGS)

software development for easy and comprehensive data analysis

Interdisciplinary group

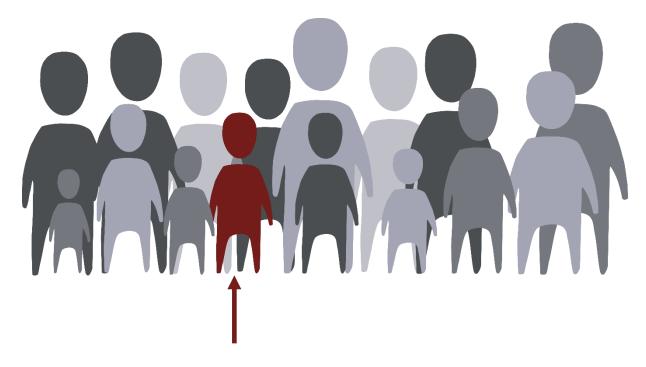
- close collaboration between clinic and research
- clinicians, biochemists, biologists, bioinformaticians

Follow-up functional studies in the wet-lab









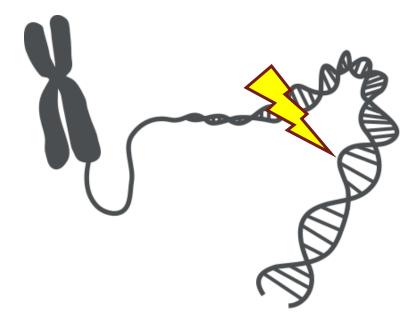
Less than 1 affected in 2.000 individuals



But: 7000 – 8000 rare diseases are known.

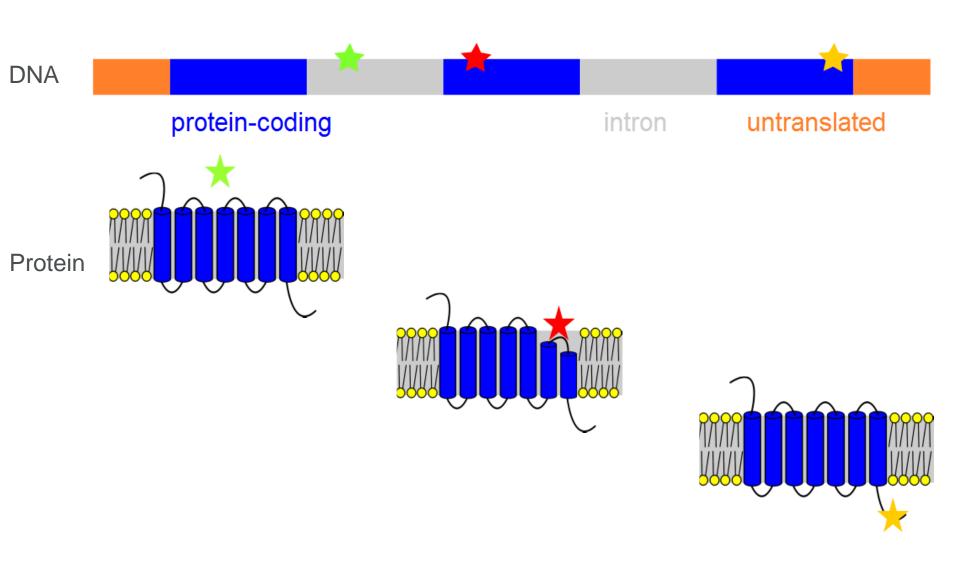
>> World-wide more than **300 Millions** patients

80% are genetically caused.



Mutations in the genome lead to a disease.

Rare diseases – genetically caused.





Severe symptoms occur congenital or in early childhood.

Diagnostics of rare diseases are often difficult and time-intensive



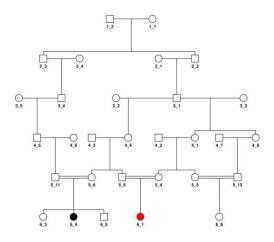
Many patients don't get a clear molecular diagnosis or no diagnosis at all.

The case of Family D.

- huge, consanguinous family from libanon
- two affected cousins
- development initially normal
- gait ataxia, distal muscle atrophy and developmental speech delays at the age of 3-5 yrs
- neurography: combined demyelinating / axonal neuropathy
- progredient course
- both patients meanwhile died



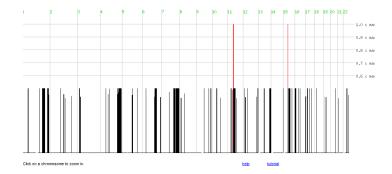




Strategies to find the disease mutation

Homozygosity mapping

Genome-wide homozygosity in Jana_D_4_samples_known_Vars - basic_copy

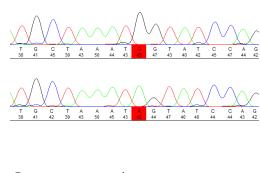


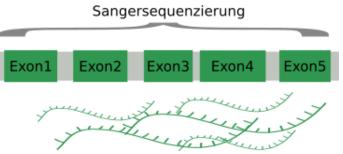
score	cm	from (pp)	to (pp)	Irom SNP	to SNP	build 37	
bros	ad - I	use this wh	en you exp	ect some ge	enetic he	terogeneit	y
240	11	48999527	51483353	n/a	n/a	region	genotypes
240	15	41829230	42235173	n/a	n/a	region	genotypes
nan	row -	use this w	hen all pati	ents are in i	the same	family	
240	11	48999527	51483353	n/a	n/a	region	genotypes
240	15	41829230	42235173	n/a	n/a	region	genotypes

Strategies to find the disease mutation

Homozygosity mapping

Sequencing of single candidate genes





Questions in 2009

How to sequence many genes at a time?

Next Generation Sequencing (NGS)

- massive parallel sequencing
- detection of tens of thousands of variants



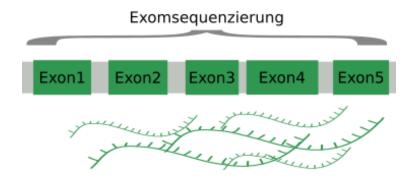
Strategies to find the disease mutation

Homozygosity mapping

Sequencing of single candidate genes

Whole Exome Sequencing





Questions in 2009

How to sequence many genes at a time?

 How to automatically evaluate the disease potential of found DNA variants?



http://www.mutationtaster.org/

Schwarz JM, Roedelsperger C, Schuelke M, Seelow D: MutationTaster evaluates disease-causing potential of sequence alterations. Nat Methods. 2010 Aug;7(8):575-6. doi: 10.1038/nmeth0810-575.













mutation t@sting

NEWS

documentation | FAQs

single query

query chromosomal positions

QueryEngine

• other applications | team

Gene	SOD1	HGNC gene symbol, NCB	l Gene ID, Ensembl gene ID <u>sho</u>	ow available tr	anscripts clear input
Transcript	ENST00000270142 Choose the transcript: ○ ENST00000470944 (, ○ ENST00000270142 (, ○ ENST00000389995 (, ○ ENST00000476106 (,	protein_coding, 966 protein_coding, 865	bases) <u>NM_000454</u> bases)		
Position / snippet refers to	coding sequence (OF	२F) 🔘 transcrip	t (cDNA sequence)	O gene	(genomic sequence)
Alteration	all types by sequence C[A/G]TGTTCATGAGT enter a few bases are Format: ACTGTC[A/T] GTGT ACTGTC[ACGT/-] GTGACTGTC[ACGT/-] GTGACTGTC[ACGT/-] GTGACTGTC[-/AA] GTGTS single base exchange to enter position	ound your alteration F GTF TGTF TF	CAGGCTGT A substituted by <i>T</i> AG substituted by <i>T</i> ACGT deleted AA inserted		options ☐ show nucleotide alignment
	and new base			J	
	insertion or deletion by position				
	enter positions of last wild type bas first wild type bas and the inserted ba	se after alteration	f applicable)		
Name of alternation					
Name of alteration	SOD1_ALS	if you would like to have	a name for this alteration in the	e output later (on, please type in here

continue







documentation







mutation tosting

Alteration SOD1_ALS

Prediction disease causing

Model: simple aae, prob: 0.99999999999447 (classification due to ClinVar, real probability is shown anyway) (explain)

Summary

HGNC symbol

DNA changes

AA changes

regulatory features

phyloP / phastCons

- amino acid sequence changed
- known disease mutation at this position (HGMD CM930680)
- known disease mutation: rs121912443 (pathogenic)
- · protein features (might be) affected
- · splice site changes

<u>analysed issue</u>	<u>anaiysis result</u>
name of alteration	SOD1 ALS

chr21:33036170A>G alteration (phys. location)

SOD1

Ensembl transcript ID ENST00000270142

UniProt peptide P00441

alteration type single base exchange

CDS alteration region

> c.140A>G cDNA.288A>G g.4236A>G

H47R Score: 0.79 explain score(s)

position(s) of altered AA

47 if AA alteration in CDS

frameshift

known disease mutation: rs121912443 (pathogenic for Amyotrophic lateral sclerosis type 1) dbSNP_NCBI variation viewer dbSNP / TGP / HGMD(public) / ClinVar

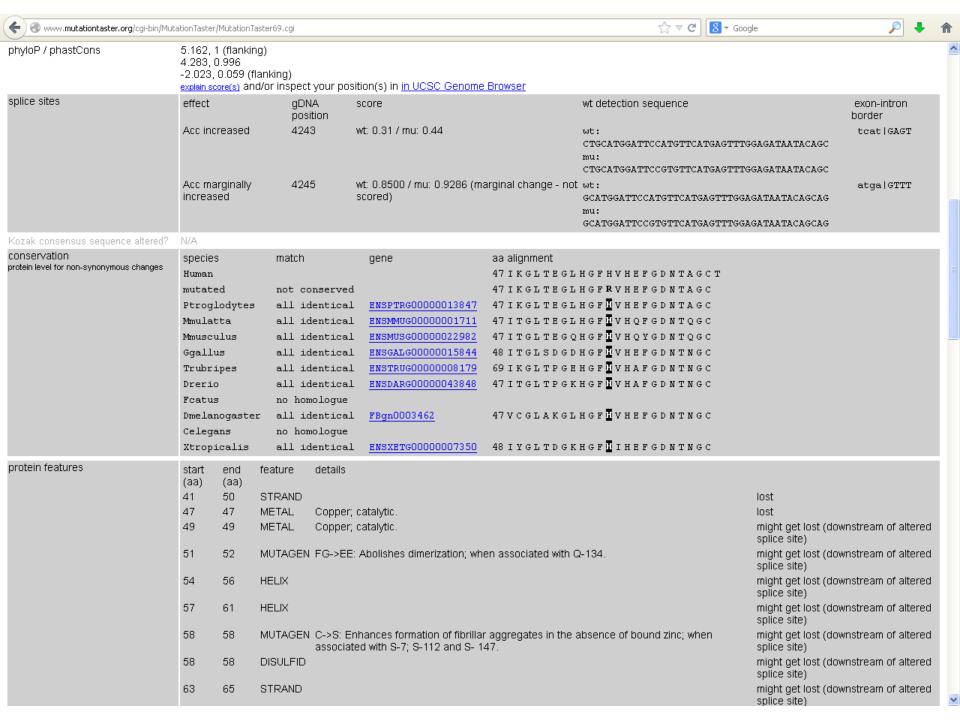
known disease mutation at this position, please check HGMD for details (HGMD ID CM930680)

H3K9me1, Histone, Histone 3 Lysine 9 mono-methylation H3K36me3, Histone, Histone 3 Lysine 36 Tri-Methylation

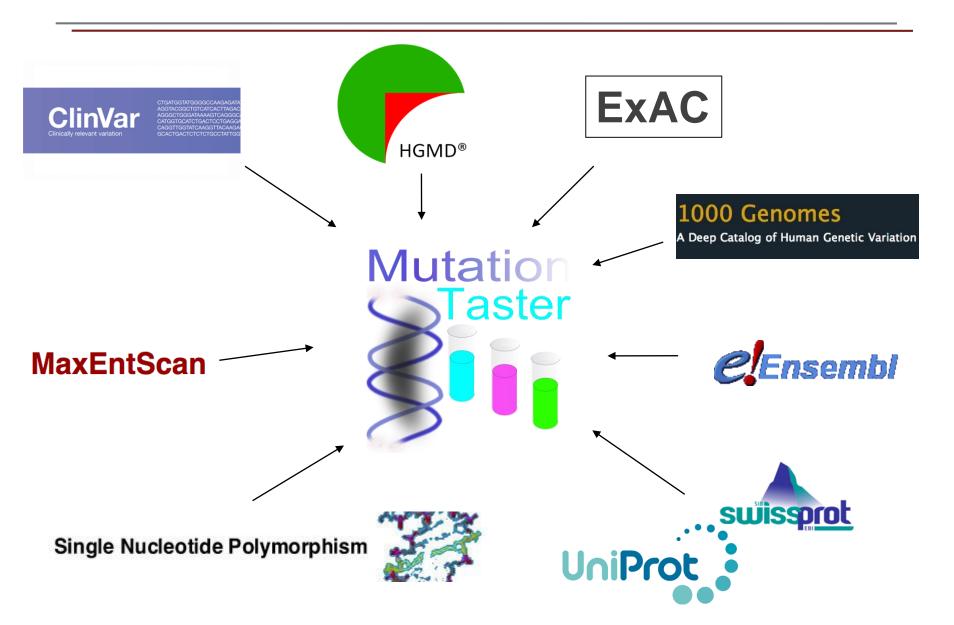
H3K79me2, Histone, Histone 3 Lysine 79 di-methylation H4K20me1, Histone, Histone 4 Lysine 20 mono-methylation

5.162, 1 (flanking) 4.283.0.996

2.022.0.050 (flanking)



Integrated data (examples)



General tests – entire transcript

test	data source / software
conservation (DNA)	BLASTN / phyloP / phastCons
splicing	MaxEntScan
occurence in data bases - known as disease mutation?	NCBI ClinVar / HGMD
occurence in data bases - known as polymorphism?	dbSNP / 1000G / ExAC
regulatory elements	Ensembl Regulation

HGMD = Human Gene Mutation Database; 1000G = 1000 Genomes Project

Tests for coding sequence (CDS)

test	data / software
amino acid exchange	MutationTaster
conservation (protein)	BLASTP
protein features	UniProtKB / SwissProt
frameshift	MutationTaster
length of protein	MutationTaster

Tests for untranslated regions (UTR)

region	test	data source / software
5'UTR	Kozak consensus sequence	MutationTaster
3'UTR	Poly(A)-Signal	polyadq













mutation tosting

Why

disease causing?

documentation

Alterat

Predic

Summary

analysed i

name of alt alteration (

HGNC symt Ensembl tra

UniProt per

alteration to

alteration re DNA chang

AA change:

position(s) if AA alteration

frameshift

dbSNP / Td ClinVar

regulatory features

H3K9me1, Histone, Histone 3 Lysine 9 mono-methylation

H3K36me3, Histone, Histone 3 Lysine 36 Tri-Methylation H3K79me2, Histone, Histone 3 Lysine 79 di-methylation

H4K20me1, Histone, Histone 4 Lysine 20 mono-methylation

phyloP / phastCons

5.162, 1 (flanking) 4.283.0.996 2.022.0.050 (flanking)









Computer based classification

- "either-or" decision
- automated classification by Bayes classifier
- simple handling
- CPU/RAM friendly
- fast, robust results

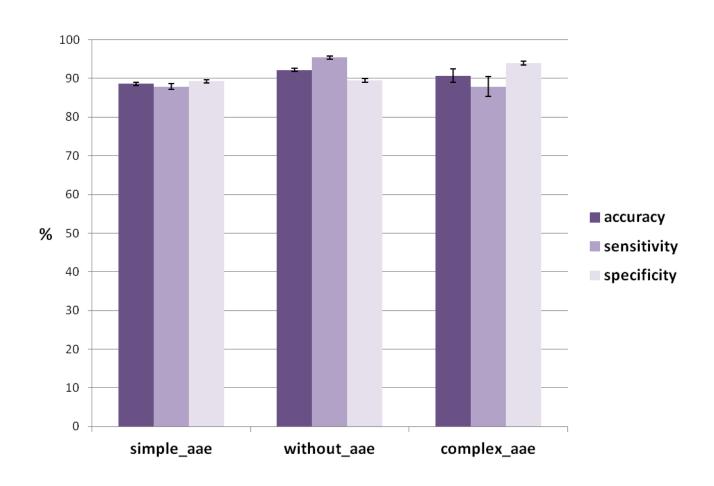
Bayes classifier in MutationTaster

- training with
 - a) > 100,000 known disease mutations

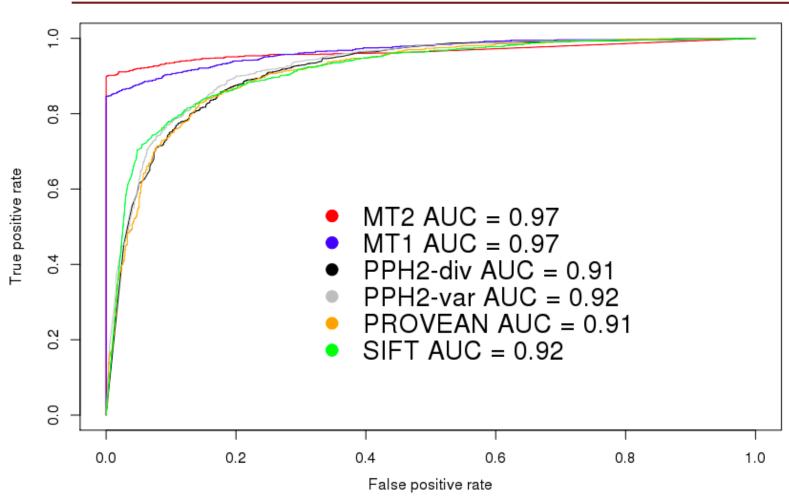
 Human Gene Mutation Database (HGMD)
 - b) > 6,000,000 harmless polymorphisms 1000 Genomes Project (1000G)

 classification model: stores frequencies of different values / test results

5-fold cross validation: results



Comparison with similar programs



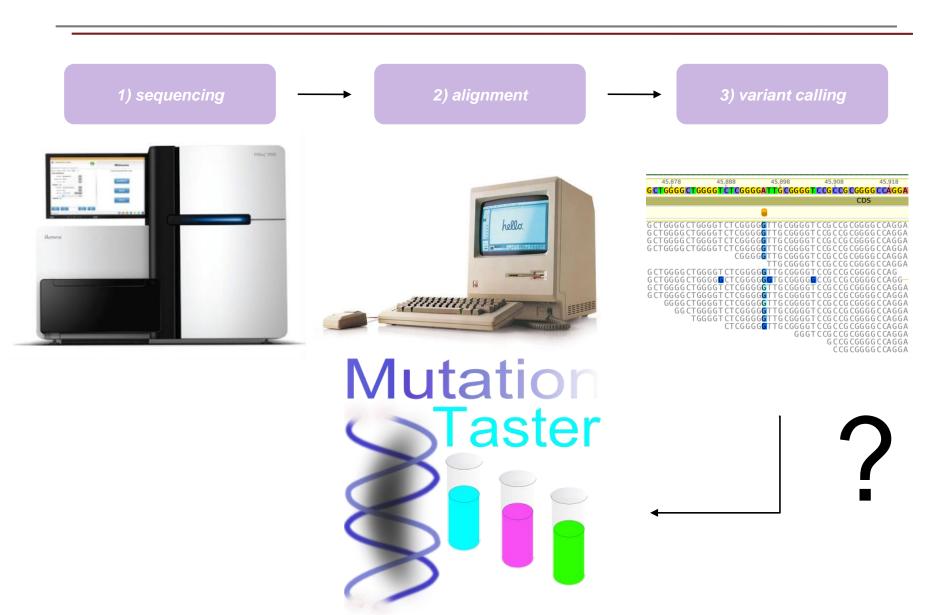
test set: 1100 polymorphisms from 1000 Genomes Project

1100 disease mutations from *Human Gene Mutation Database* (HGMD)

abbreviations: *PPH2-var* = PolyPhen2 HumVar model, *PPH2-div* = PolyPhen2 HumDiv model

MutationTaster's, automatic classifications based on public variant data from 1000G and ClinVar were ignored for this comparison.

NGS data and MutationTaster



MutationTaster QueryEngine



minimum coverage

Submit

mutation t@sting

the results (archived as .zip) will be send via E-mail to you. For this reason, you have to provide a valid E-mail address. Look up more details in the documentation.

QueryEngine

- documentation | FAQs
- single query
- · query chromosomal positions
- QueryEngir
- . MutationDistiller (public beta)
- . Regulation Spotter (public beta)
- other applications | team
- slides ESHG2017 Copenhagen

VCF file Durchsuchen... MTQE_example_vcf.vcf clear input Please zip or gzip large files! Format: #CHROM ID REF ALT QUAL FILTER INFO $\texttt{DP=2;AF1=0.5003;CI95=0.25,0.75;DP4=1,0,0,1;MQ=60;FQ=-3.1;PV4=1,1,1,1} \\ \texttt{GT:PL:DP:GQ} \\ \texttt{0/1:33,0,28:2:30} \\ \texttt{0/1:33,0,28:2:3$ (tab delimited) The coordinates must refer to GRCh37 (also called hg19). Project name Example queue status 2017-09-21 17:42 CEST (refresh) low load - our bored server is looking forward to Unfortunately, E-Mail notification does currently not work. We are sorry for the inconvenience and try to solve the problem as soon as analysing your data possible 0 jobs running, 0 queued, 1511.2 millions alterations analysed. free slots (0 running, 16 of 16 available - at 0% capacity) large jobs free slots (0 running, 20 of 20 available - at 0% capacity) E-mail address free slots (0 running, 30 of 30 available - at 0% capacity) small jobs generate HTML files (slower) Analysis settings search for homozygous variants yes analyse complete VCF ...but only exons with 10 bases intron flanking analyse variants on chr combine neighbouring variants yes ...but only exons with 10 bases intron flanking filter against TGP analyse custom regions (select to enter) or more TGP samples homozygous in 4 exclude custom regions (select to enter) or more TGP samples heterozygous in 20

We offer automated MutationTaster analysis of variants from Next Generation Sequencing projects. Variants must be in VCF format and refer to GRCh37 / hg19. After your VCF file has been analysed, the link to download



[MutationTaster] for Example - explore your results

highly-parallel fashion (500,000 variants / hour)

Statistics	Display / filter / export res	Query GeneDistiller	
submitted variants 4 pre-discarded variants 0 analysable alterations 4 discarded (TGP) 1 discarded (out of gene/exon/region) 0 alterations analysed 3 MT cases 9 type without_aae 6 type simple_aae 3	sort & group sort & group sort & group by prediction model gene symbol sort & group by prediction model gene symbol variation sort by these attributes 1: descending 2: descending 3: descending	hide silent alterations all predicted polymorphisms known polymorphisms prediction problems show only homozygous alterations all recessive (homozygous, comp-het, genes with >1 disease-causing variant) only genes relevant for mitochondria	call GeneDistiller for all disease-causing mutations non-synonymous disease-causing mutations non-synonymous alterations (excl. known polymorphisms)
prediction disease_causing_automatic 4 prediction polymorphism 5 genes hit (except HapMap/TGP)non-synonymous 2disease causing 2disease causing, non-synonymous 2 Analysis options	display (5000 results per page)	get the data export as TSV	
coverage threshold 4 filter out variants homozygous in TGP 4 >			

- . see QueryEngine log for further details on the analysis run
- · were discarded before mapping
- · were filtered out (TGP)

- · were not mapped to a transcript
- · got no MutationTaster prediction



http://www.mutationtaster.org/

Schwarz JM, Cooper DN, Schuelke M, Seelow D:

MutationTaster2: mutation prediction for the deep-sequencing age.

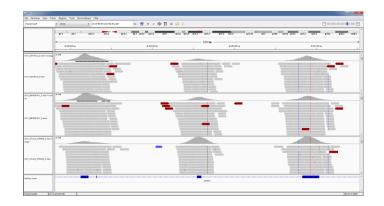
Nature Methods (April 2014) 11(4):361-2.

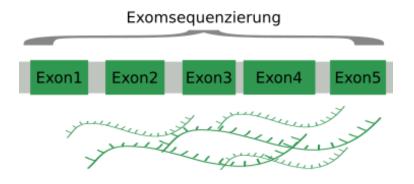
Strategies to find the disease mutation

Homozygosity mapping

Sequencing of single candidate genes

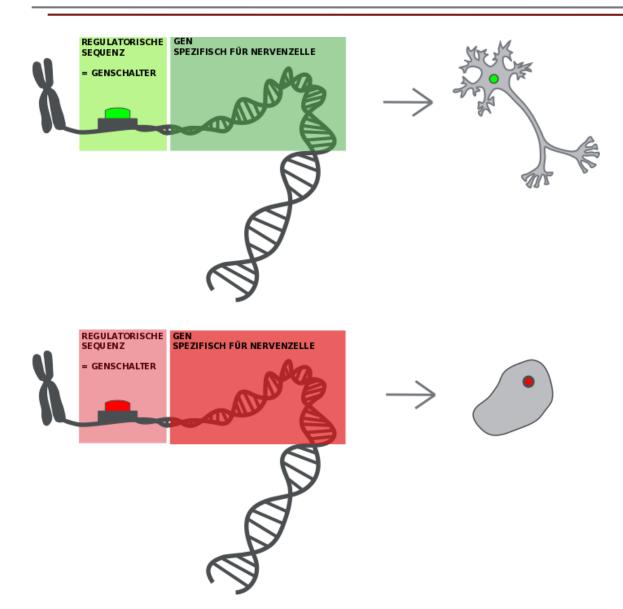
Whole Exome Sequencing



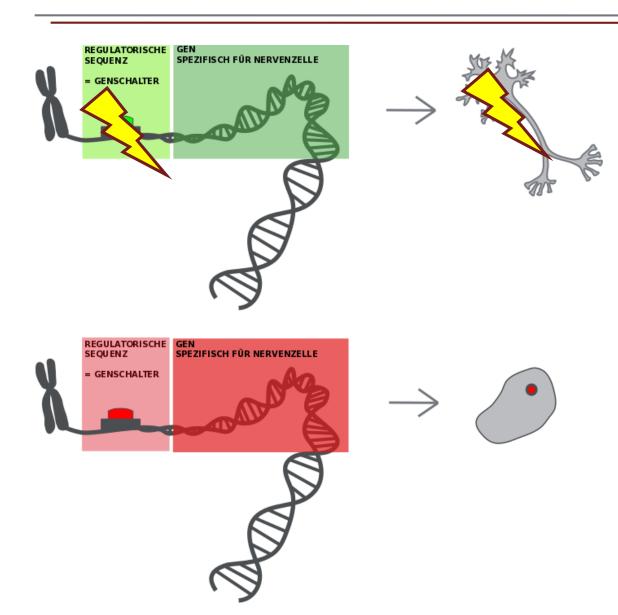


....did not detect the disease mutation!

Extragenic / regulatory DNA variants

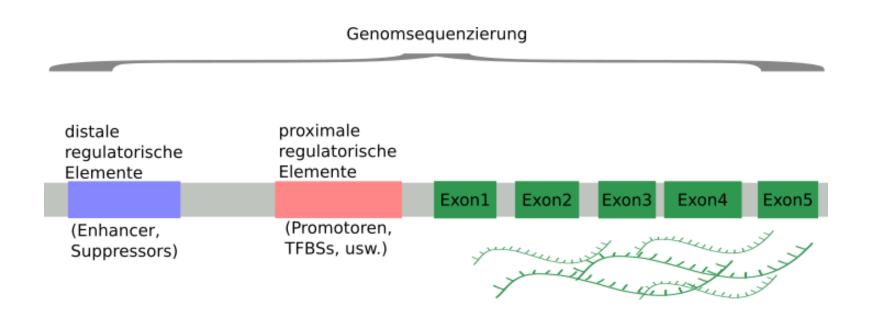


Extragenic / regulatory DNA variants



Mutations in regulatory elements can lead to onset of a genetic diseases.

Whole Genome Sequencing



Strategies to find disease mutations

Homozygosity mapping

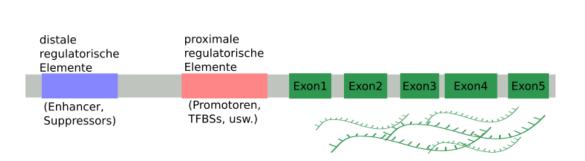
Sequencing of single candidate genes

| Companies | Table | Suprementary | Table |

Whole Exome Sequencing

Whole Genome Sequencing

Genomsequenzierung



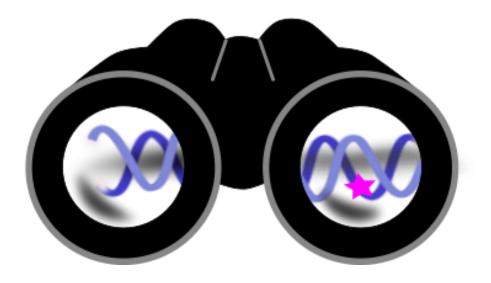
~ 3 million variants

Mutation

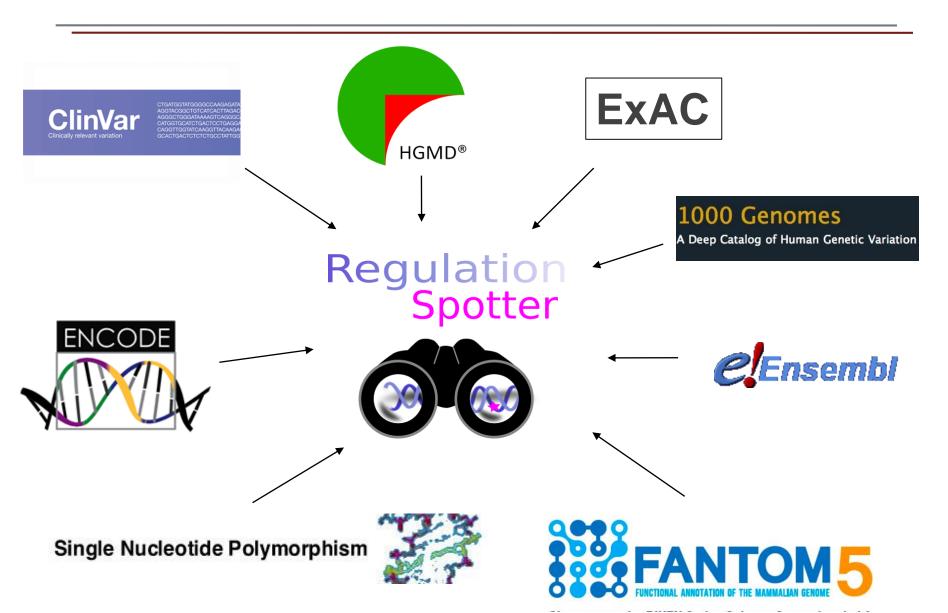
Only for genes.

Schwarz JM, Cooper DN, Schuelke M, Seelow D: MutationTaster2: mutation prediction for the deep-sequencing age. *Nature Methods* (April 2014) 11(4):361-2.

Regulation Spotter



Integrated data (examples)



Integrated data (examples)

Integrated regulatory features

- 125 regulatory features including
 - promoter and enhancer annotations
 - histone modifications
 - DNAsel hypersensitive sites
 - Transcription factor binding sites

X-score

- reflects amount of evidence that a variant is located in a regulatory relevant region
- calculated based on 54 most meaningful regulatory annotations / features

feature weight:

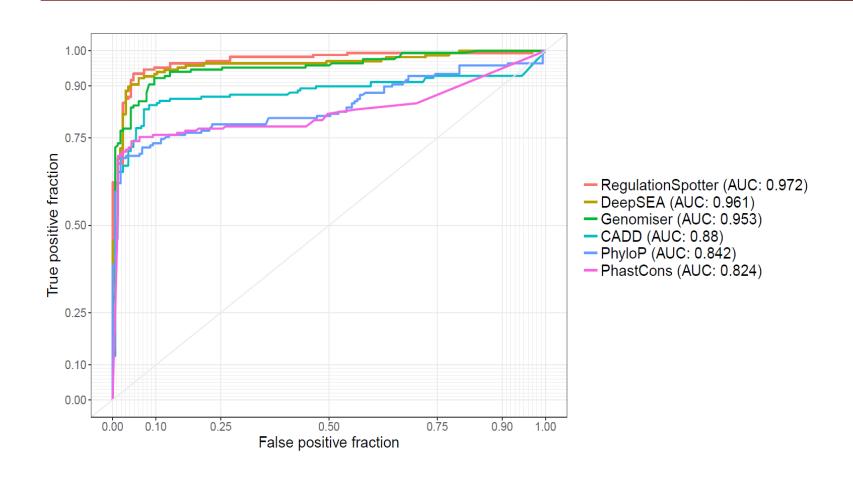
- reflects assumed impact of each feature
- determined based on relative risk to occur in a testset of functional vs. non-functional extragenic variants
- X-score represents the sum of the weights for all features annotated for a variant

X-score

- low number of extragenic variants validated as regulatory disease mutations
 - not enough for classifier training

positive: 295 disease variants from HGMDpro

Comparison with similar programs

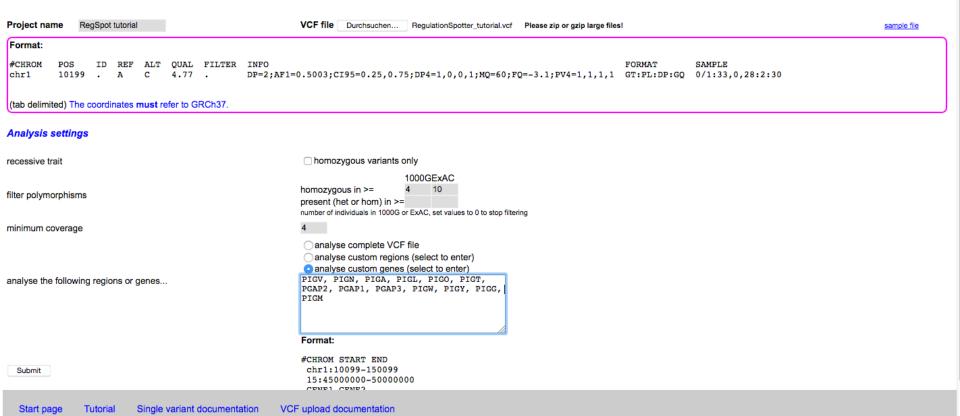


test set: 166 extragenic harmless polymorphisms from *1000 Genomes Project* 166 extragenic disease mutations from NCBI *ClinVar*

RegulationSpotter – submit VCF file



We offer automated analysis of variants from Whole Genome Sequencing projects in either user-defined genomic regions or in candidate genes and the genomic elements interacting with these. Variants must be in VCF format and refer to GRCh37. Look up more details in the documentation. You can also follow the provided tutorial to become familiar with the software and the results.



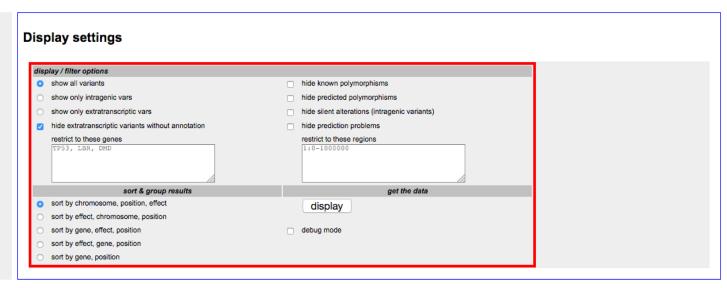
RegulationSpotter – results synopsis

Regulation Spotter

[RegulationSpotter] - synopsis & display settings

RegulationSpotter documentation





see QueryEngine log for further details on the analysis run

· were not mapped to a transcript

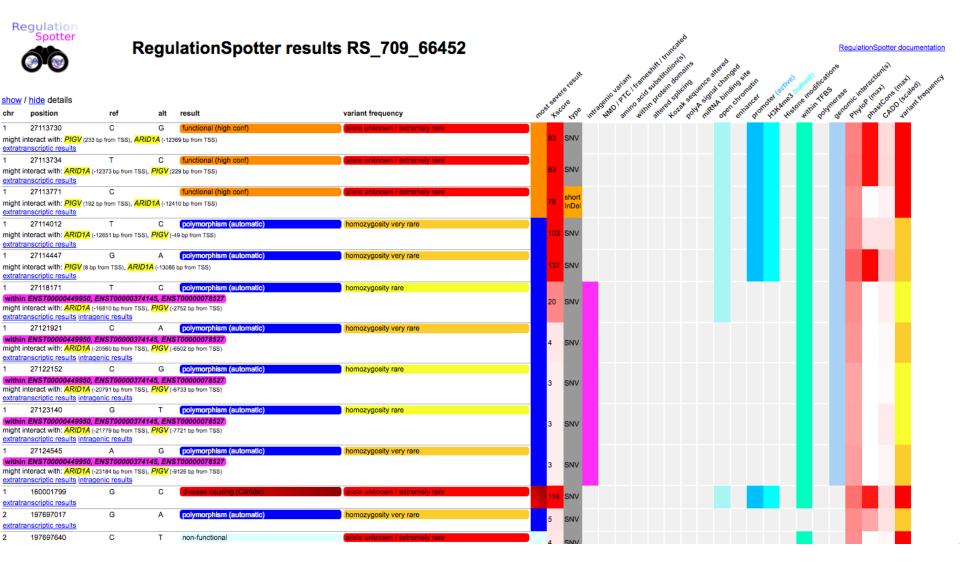
· were discarded before mapping

project ID 709_66452

elete your project

export your project

RegulationSpotter – results matrix



RegulationSpotter – detailed results



histone modifications

regul@tion spotting

5 cell types: CM12878 USMM USMMtube K562 Octoobl

functional (high confidence) Likely effect Model: extratranscriptic, score: 117.35756097561 Summary hyperlink within TFBS · within active promoter analysed issue analysis result name of alteration alteration (phys. location) chr17:37844280G>C IGV alteration type single base exchange alteration region extratranscriptic known variant Reference ID: rs564810973 database homozygous (C/C) heterozygous allele carriers 0 1000G ExAC promoters Ensembl Promoter ENSR00000126595 ctive A549, Aorta, B cells (PB) Roadmap, CD14+CD16- monocyte (CB), CD14+CD16- monocyte (VB), CD4+ ab T cell (VB), CD8+ ab T cell (CB), CM CD4+ ab T cell (VB), DND-41, eosinophil (VB), EPC (VB), erythroblast (CB), Fetal Adrenal Gland, Fetal Intestine Large, Fetal Intestine Small, Fetal Muscle Leg, Fetal Muscle Trunk, Fetal Stomach, Gastric, GM12878, H1ESC, H1-mesenchymal, H1-neuronal progenitor, H9, HeLa-S3, HepG2, HMEC, HSMM, HSMMtube, HUVEC, HUVEC prol (CB), IMR90, iPS-20b, iPS DF 19.11, iPS DF 6.9, K562, Left Ventricle, Lung, M0 macrophage (CB), M0 macrophage (VB), M1 macrophage (CB), M1 macrophage (VB), M2 macrophage (VB), M3 macrop myelocyte (BM), NH-A, NHDF-AD, NHEK, NHLF, Osteobl, Ovary, Pancreas, Placenta, Psoas Muscle, Right Atrium, Small Intestine, Spleen, T cells (PB) Roadmap poised Fetal Thymus, H1-trophoblast, neutrophil (CB), neutrophil (VB), Thymus FANTOM predictions FANTOM TSS relaxed 17:37844268-37844316(-) RegulationSpotter 17: 37843854-37844404 H3K4me3 ERBB2 DNase1 ENST00000584014 17: 37843893-37844443 DNase1 H3K4me3 ERBB2 ENST00000406381 DNase1 H3K4me3 ERBB2 ENST00000578199 17: 37843872-37844422 17: 37844252-37844802 DNase1 H3K4me3 PGAP3 ENST00000582276 17: 37844204-37844754 DNase1 H3K4me3 PGAP3 ENST00000584620 17: 37844260-37844810 DNase1 H3K4me3 PGAP3 ENST00000579146, ENST00000429199, ENST00000378011, ENST00000309862 enhancers none found epigenetic marks DNase1 DNase1 hypersensitive site in promoter 17:37843652-37844692 (RegulationSpotter) H3K4me3 H3K4me3 in 3 cell lines, overlap with promoter 17:37843348-37845639

Summary

MutationTaster

- evaluates disease potential of intragenic DNA variants
- intronic & exonic
- non-synonymous / synonymous and silent single base exchanges and indels
- classification by Bayes Classifier
- for single variants or complete vcf files
- for variants involved in Mendelian disorders!

http://www.mutationtaster.org/

RegulationSpotter

- evaluates functional / regulatory impact of extragenic DNA variants
- single base exchanges and indels
- classification with help of X-score
- for single variants or complete vcf files

http://www.regulationspotter.org/

Thank you!

AG Seelow (Charité Berlin / BIH)

Daniela Hombach Sebastian Köhler Dominik Seelow

https://translationalgenomics.charite.de/

AG Schülke (Charité Berlin / Neuropeadiatrics)

Ellen Knierim Gudrun Schottmann Markus Schülke

HGMD

David Cooper (University of Cardiff)

Practicals:

Wednesday

11.45-13.15

14.15-15.45

Zürich 3

Schwarz JM, Rödelsperger C, Schuelke M, Seelow D.

MutationTaster evaluates disease-causing potential of sequence alterations.

Nature Methods. 2010 Aug;7(8):575-6

Schwarz JM, Cooper DN, Schuelke M, Seelow D.

MutationTaster2: mutation prediction for the deep-sequencing age.

Nature Methods. 2014, Apr;11(4):361-2.

Schwarz JM, Hombach D, Koehler S, Cooper DN, Schuelke M, Seelow D.

RegulationSpotter: Easy and comprehensive analysis of extratranscriptic sequence variants.

Under revision.

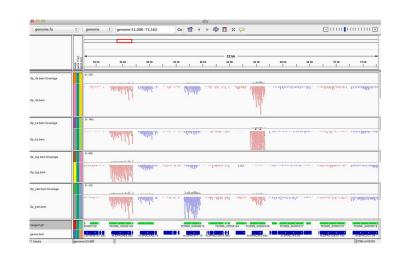
The case of Family D. – what comes up next

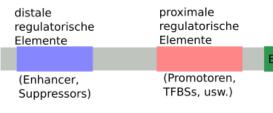
Homozygosity mapping

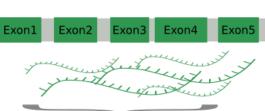
Sequencing of single candidate genes

Whole Exome Sequencing

Whole Genome Sequencing







Transkriptomsequenzierung

iPSCs

Whole Transcriptome Sequencing