



VAAST

Variant Annotation Analysis and Search Tool

CSHL Programming for Biology Oct 2014

Barry Moore
Director, Science & Research
USTAR Center for Genetic Discovery
Department of Human Genetics
University of Utah



Outline

- Historical Perspective
- Variant Calling (Follow the Probabilities)
- VAAST
 - VAT
 - VST
 - VAAST 2.0
- Rare Disease Applications
- Common Disease Applications
- Future Directions

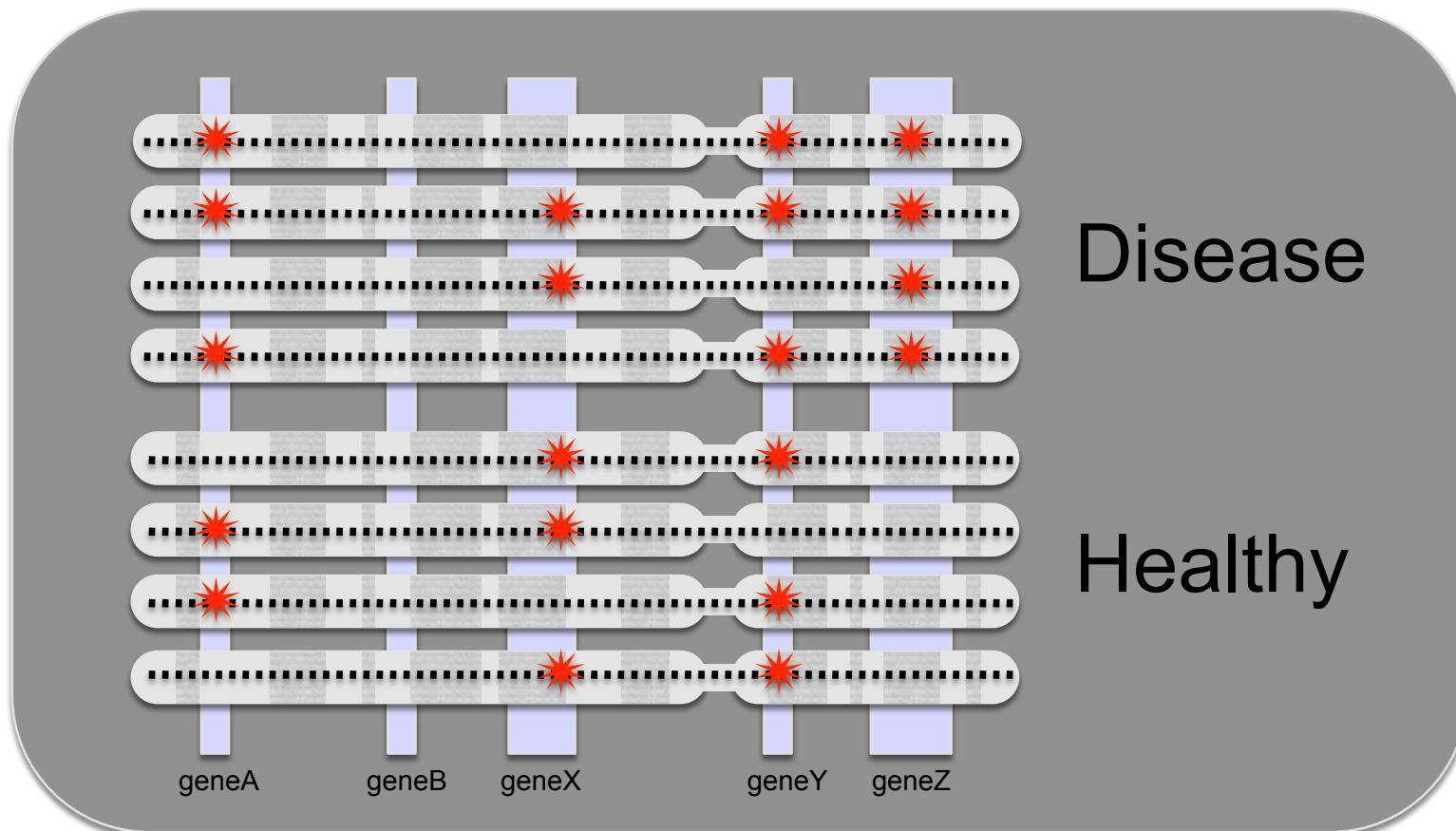
Motivation

- Billions of years of evolution have fine tuned our DNA sequence.
- Genetic alterations to that sequence can cause disease.
- Knowing which mutations provides clues to understanding the disease.
- What are the mutations – developing technology.
- Which mutation – developing analysis methodologies.

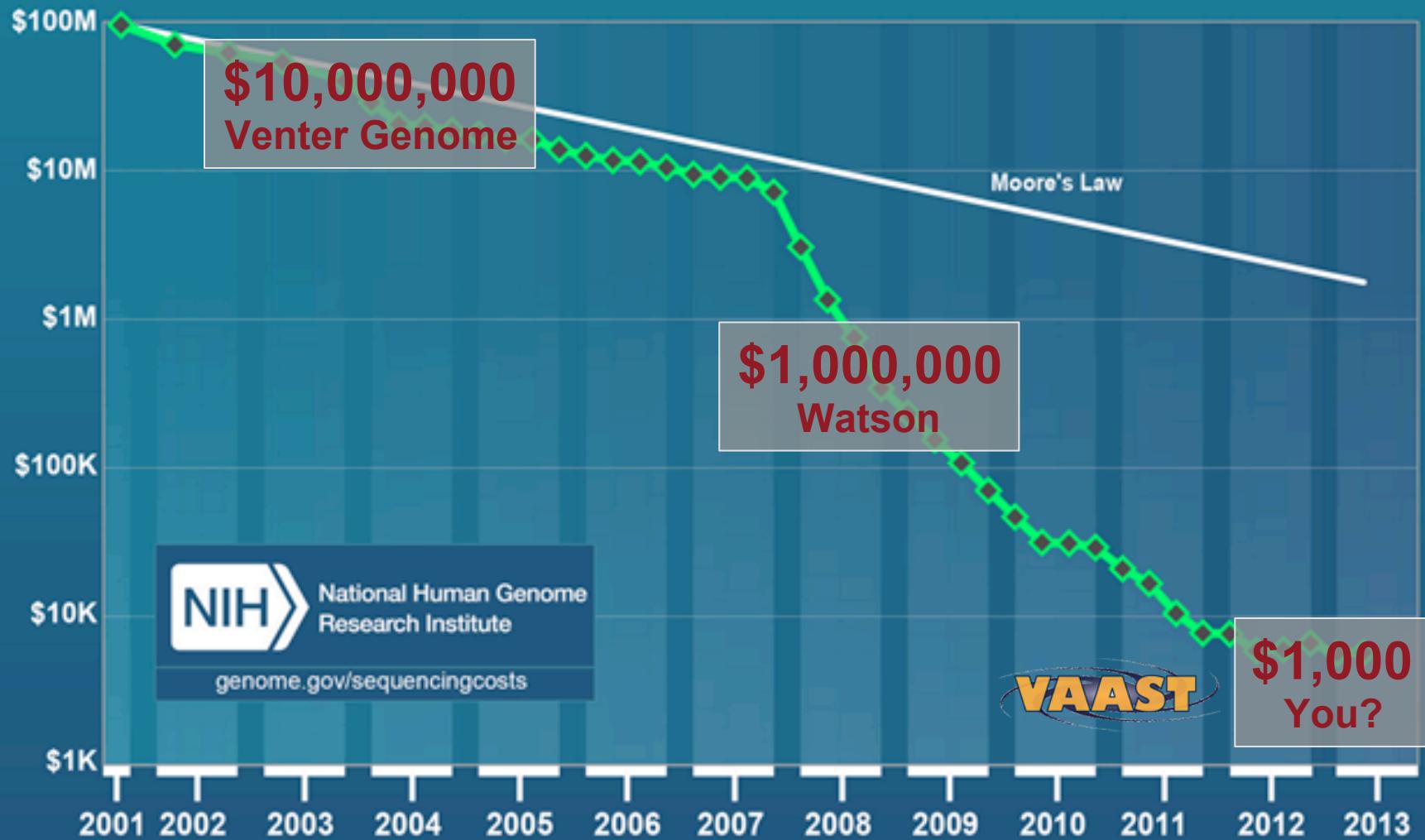
A Very Breif History of Medical Genetics

- Cytogenetics
 - Downs Syndrome - 1959
- Linkage Mapping
 - HTT (Huntington's Disease) gene mapped – 1983
- Positional Cloning
 - CFTR (Cystic Fibrosis) gene discovered - 1989
- Sequencing, microarrays and GWAS
 - ARMD, CD, MI, IBD – 2005-2006
- Next generation sequencing – personalized genomics
 - Charcot-Marie-Tooth, Miller Syndrome, Ogden Syndrome 2010-2011

Genome Wide Association



Cost per Genome



VAAST Overview

- Probabilistic tool for disease gene discovery
- Aggregative variant analysis – feature based
- Both allele and AAAS frequencies
- Conservation-controlled AAS
- Implements numerous filters
- Standardized ontology based formats
- Modular and flexible in design

Outline

- Historical Perspective
- Variant Calling (Follow the Probabilities)
- VAAST
 - VAT
 - VST
 - VAAST 2.0
- Rare Disease Applications
- Common Disease Applications
- Future Directions

NGS Data Analysis

- Trillions of glowing DNA fragments produce base calls (reads).
- 100s of millions of sequence reads produce alignments.
- 10s of millions of variant sites produce variant calls.
- 10s of thousands of variants analyzed for association with disease.
- 1 gene causes disease.

Follow the Probabilities

- Base calling – Bayesian inference
- Base quality score recalibration – Covariant analysis
- Mapping quality – Pseudo-probabilistic
- Variant calling – Bayesian inference
- Variant quality score recalibration – LOD ratio based on a trained Gaussian mixture model
- VAAST - CLRT

Variant calling - Individual

- What is the probability that this site is reference vs. variant given the reads aligned at this site.
- What is the probability that this site has homozygous reference genotype vs. heterozygous genotype given the reads aligned at this site.

Variant calling - Population

- What is the probability that this site is reference (**FOR THE POPULATION**) given:
 - The reads aligned at this site **FOR THE POPULATION**
- What is the probability that the genotype is homozygous reference (**FOR THE INDIVIDUAL**) given:
 - The reads aligned at this site **FOR THE INDIVIDUAL**
 - The probability that this site is variant **FOR THE POPULATION**

Variant calling - Population

Reference Sequence

TACAGCTGAACTGAACTCCTGCCTGTACAGCTCGTTTCTACAAGATTCCAGACCTGGAA

Individual

TGCCTGTACAGCTCGTTTCTACAAGATTTC

AACTGAACTCCTGCCTGTAC**G**GCTCGT

TGAACCTCCTGCCTGTAC**G**GCTCGTTTCTA

TGTAC**A**GCTCGTTTCTACAAGATTCCAGA

CTCCTGCCTGTAC**G**GCTCGTTTCTACAAG

ACTCCTGCCTGTAC**G**GCTCGTTTCTACAA

GCCTGTAC**G**GCTCGTTTCTACAAGATTCC

TGCCTGTACAGCTCGTTTCTACAAGATTTC

AACTGAACTCCTGCCTGTAC**G**GCTCGT

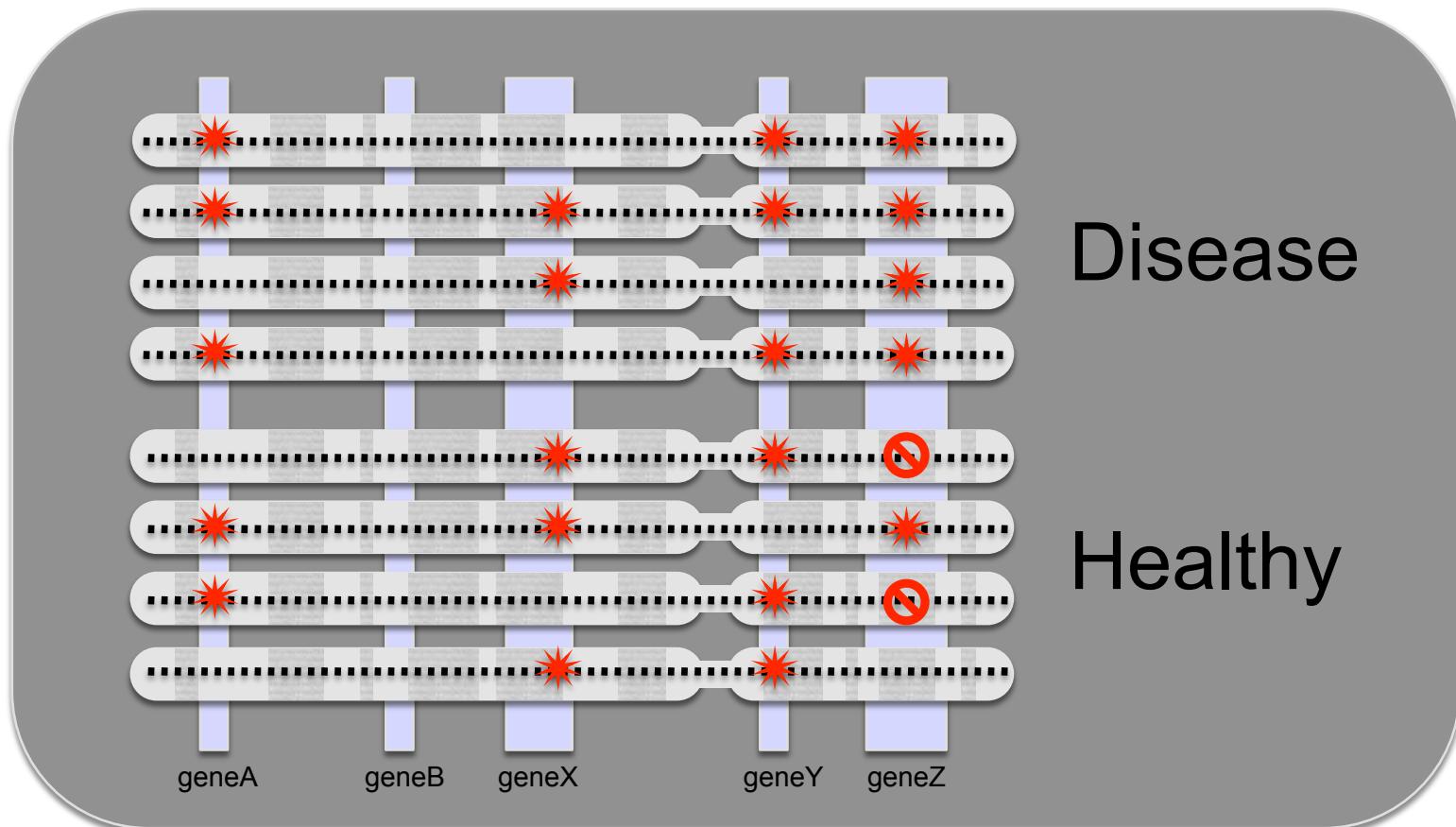
TGAACCTCCTGCCTGTACAGCTCGTTTCTA

Population

Missing data

- Low/no sequence coverage
- Low base qualities
- Variant callers typically emit only variant sites
- What happens when we don't distinguish between missing data and reference sites
- Annotating no-calls can help solve this problem
- Population based variant calling provides this

Missing data



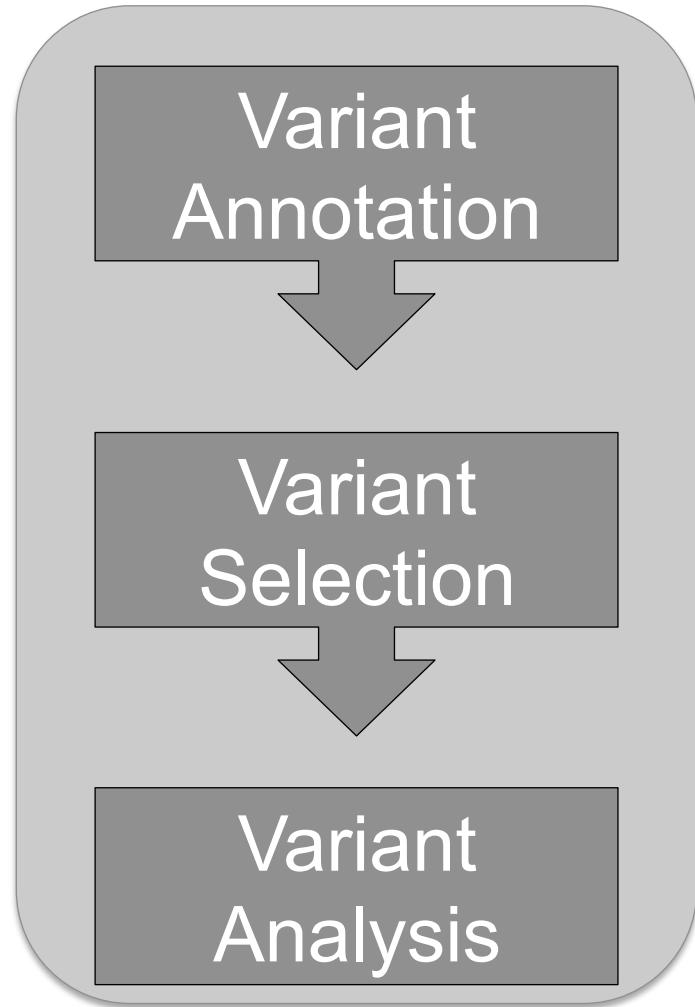
Variant Calling Format (VCF)

```
##fileformat=VCFv4.1
##FILTER=<ID=LowQual,Description="Low quality">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Approximate read depth">
##contig=<ID=1,length=249250621>
```

CHROM	POS	ID	REF	ALT	QQUAL	FILTER	INFO	FORMAT	INDIVIDUAL	INDIVIDUAL
1	745370	.	TA	T	1310.90	PASS	DP=210;	GT:GQ	0/1:99	0/1:99
1	749592	.	G	A	20	LowQual	DP=7;	GT:GQ	./.	1/1:6
1	749683	.	C	T	602.40	PASS	DP=69;	GT:GQ	./.	1/1:13
1	749856	.	C	T	261.37	PASS	DP=79;	GT:GQ	1/1:9	0/1:99
1	749899	.	G	A	302.28	PASS	DP=53;	GT:GQ	0/0:9	0/1:99
1	752566	.	G	A	1047.91	PASS	DP=47;	GT:GQ	1/1:30	0/1:29
1	752721	.	A	G	7625.90	PASS	DP=360;	GT:GQ	1/1:99	0/1:99

VAAST

- A tool for identifying disease genes and variants
- Collaboratively developed
 - Mark Yandell (University of Utah)
 - Chad Huff (MD Anderson)
 - Martin Reese (Omicia Inc.)
- Inputs
 - Target genome variant files
 - Background genome variant files
 - Genomic features (gene models)
 - Genomic sequence
- Outputs a prioritized list of features (genes/transcripts) associated with the disease genomes.
 - VAAST Score
 - P-value
 - Confidence interval

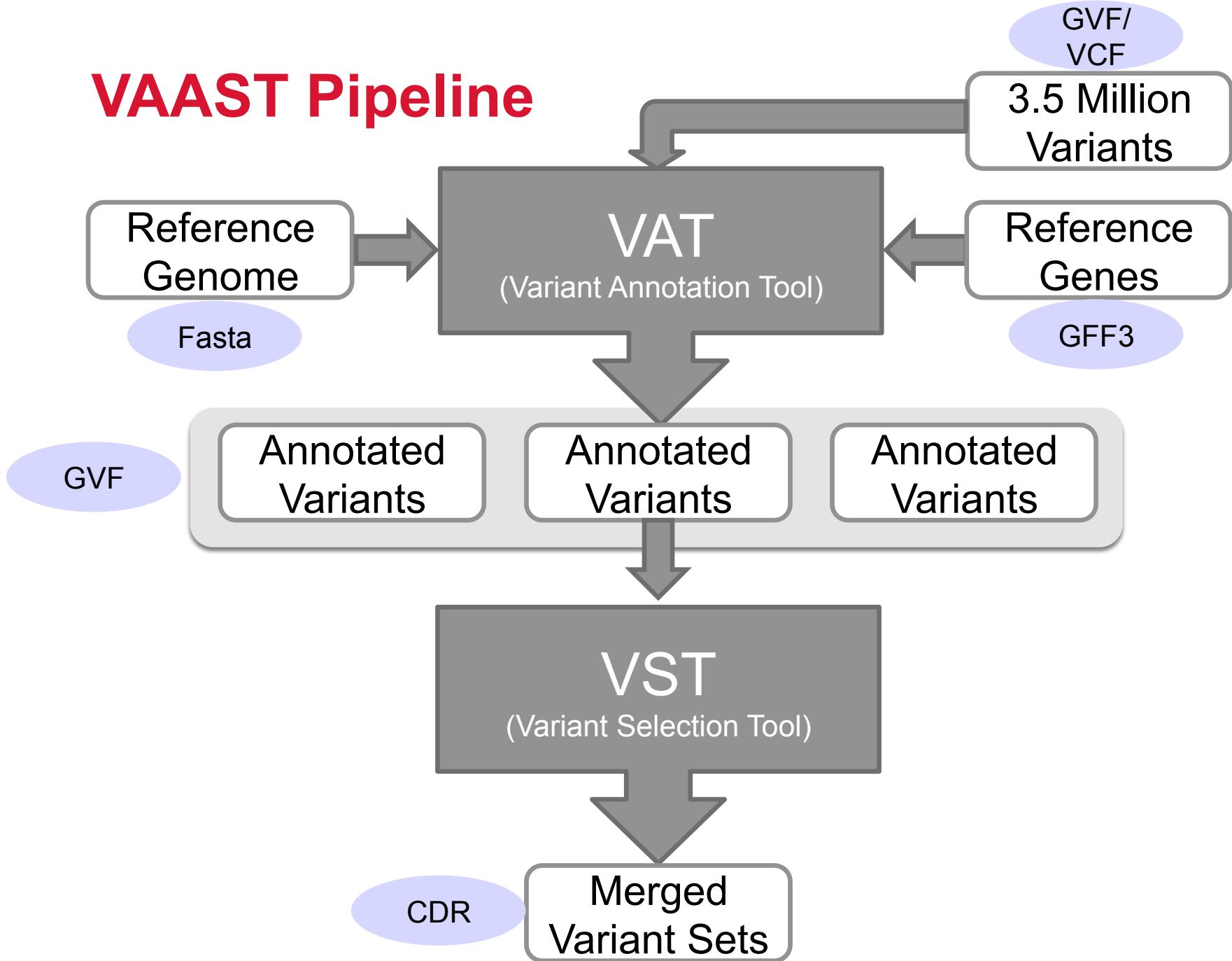


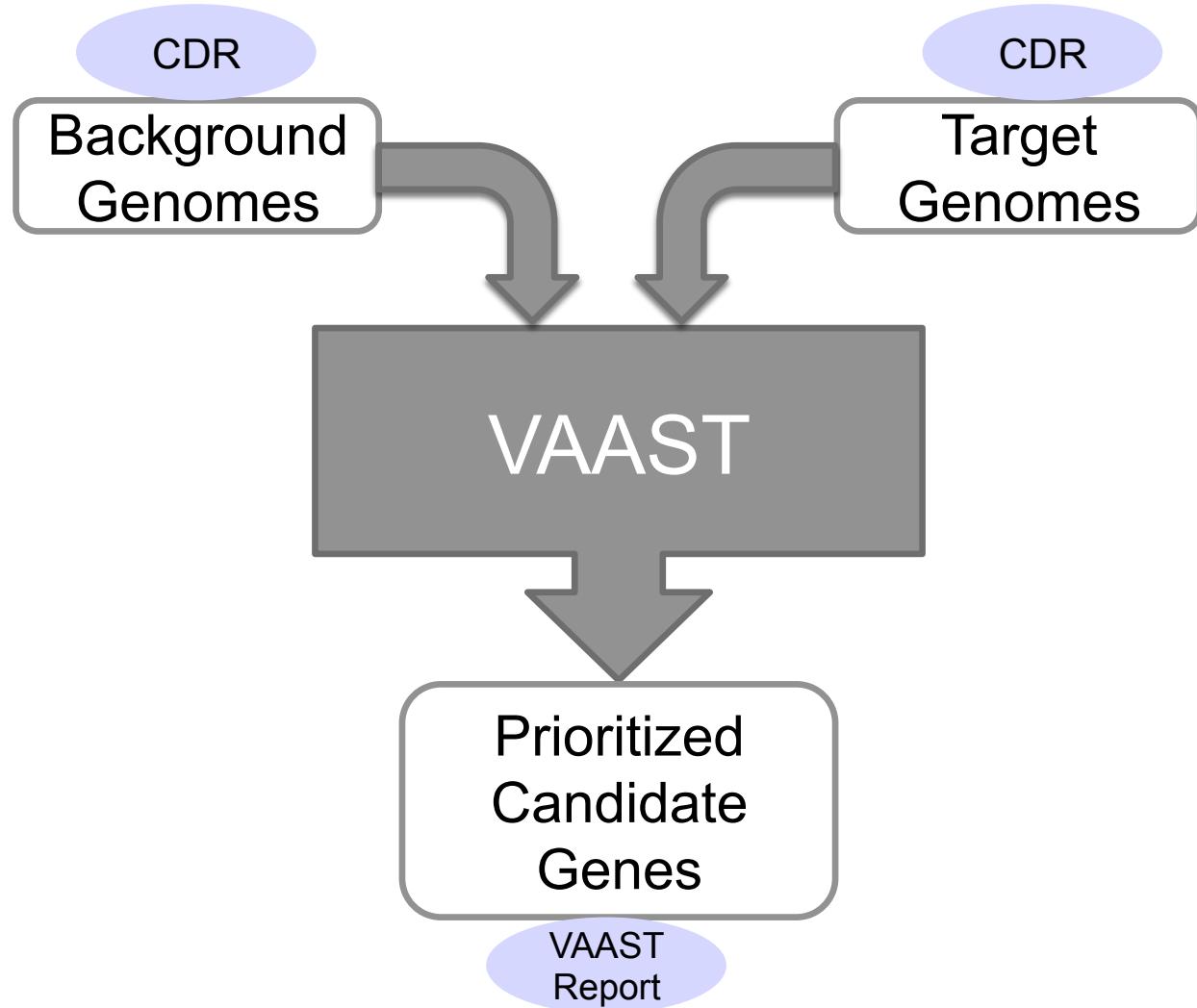
*Variant
Annotation
Tool*

*Variant
Selection
Tool*

*Variant
Annotation
Analysis
Search
Tool*

VAAST Pipeline





VAAST File Types

■ Input

- Fasta - Genome
- GFF3 - Genes
- GVF - Variants

■ Output

- GVF – Annotated Variants
- CDR – Population Variants
- VAAST – Prioritized gene list

Fasta

```
>chr1
TAACAAAATAAGATCCAGAAACTTCCATTAGCGTGGGGGTGACCATGAA
ATGCCTGGTCAAAAACCCGGGCACTGATTGTATAACCATTATGCAACTG
GTGTTGCGTCCATCAGAATCTAGTTAAGAATACTCTTCTCTATAGGA
GTCTTCGCGGCAGACCTAGCCTGCTCTGTGTCCCTGAAATGAAGGAAT
GTTCTCTCCCATTATTCTTAACAGCTTGGTTAGCAAGCTCCGCCCTC
TTCTTATCTGACCTCTAACGACCTCACCAAGATGTGTGAAGCAGCCGG
CTCCATGTGTATCAGgcacgcacgcacacacgcacacCAACCTGCA
AAGGAAATAACGGGGCAGCCCTGCAGTGTAAAAAGCAATGGGATTTGTG
GGTTCCACCTCCTCACCTAACGATCCCTGGTCTACGCTATGTCACGACCC
TCTGCTGAACCACGTCAGGGTGAACCCNNNNNNNNNNNNNNNNNNNNNNNN
```



[Home](#) [Browser](#) [Wiki](#) [GFF3](#) [GVF](#) [Resources](#) [Software](#) [About](#) [Request A Term](#) [Site Map](#)

Welcome to the Sequence Ontology

This is the home page of the Sequence Ontology (SO). SO is a collaborative ontology project for the definition of sequence features used in biological sequence annotation. SO was initially developed by the [Gene Ontology Consortium](#). Contributors to SO include the [GMOD](#) community, model organism database groups such as [WormBase](#), [FlyBase](#), [Mouse Genome Informatics](#) group, and institutes such as the Sanger Institute and the EBI. Input to SO is welcomed from the sequence annotation community. SO is also part of the Open Biomedical Ontologies library. Our aim is to develop an ontology suitable for describing the features of biological sequences. For questions, please send mail to the [SO developers mailing list](#). For new term suggestions, please use the [Term Tracker](#).

Introduction

The Sequence Ontology is a set of terms and relationships used to describe the features and attributes of biological sequence. SO includes different kinds of features which can be located on the sequence. Biological features are those which are defined by their disposition to be involved in a biological process. Examples are **binding_site** and **exon**. Biomaterial features are those which are intended for use in an experiment such as **aptamer** and **PCR_product**. There are also experimental features which are the result of an experiment. SO also provides a rich set of attributes to describe these features such as "polycistronic" and "maternally imprinted".

News

- ▶ **October 2013** GVF was used in the clinical annotation of a whole genome, for precision medicine. Integrating precision medicine in the study and clinical treatment of a severely mentally ill person

- ▶ **September 2013** The SO development team and the Monarch Initiative held a collaborative workshop in Portland to align the SO and the GENO ontologies for better annotation of phenotypes. This meeting was partially funded by the Phenotype RCN .

Variants, Features and Effects

Variant Type

- sequence_alteration
- deletion
- insertion
- duplication
- inversion
- substitution
- SNV**
- MNP
- complex substitution
- translocation

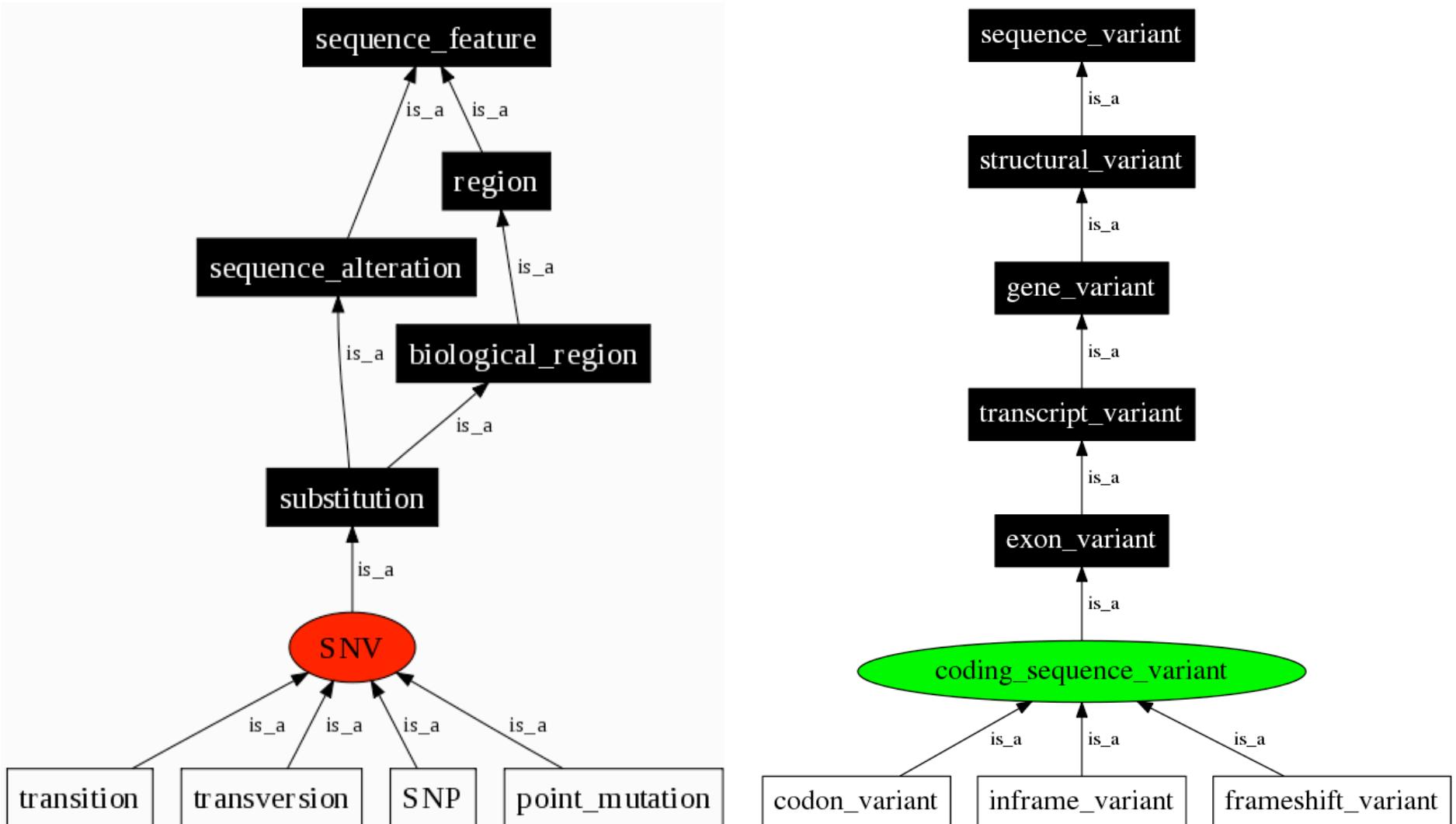
Feature Type

- sequence_feature
- gene
- mRNA
- exon
- CDS**
- splice site
- ncRNA

Variant Effect

- sequence_variant
- gene_variant
- five_prime_UTR_variant
- three_prime_UTR_variant
- exon_variant
- splice_region_variant
- splice_donor_variant
- splice_acceptor_variant
- intron_variant
- coding_sequence_variant
- stop_retained
- stop_lost
- stop_gained
- synonymous_variant
- missense_variant**
- amino_acid_substitution
- frameshift_variant
- inframe_variant

Sequence Ontology



Generic Feature Format (GFF3)

##gff-version 3								
##sequence-region chr1 1 1497228								
Seqid	Source	Type	Start	End	Score	Strand	Phase	Attributes
chr1	.	gene	1000	9000	.	+	.	ID=gene00001;Name=EDEN 3
chr1	.	mRNA	1050	9000	.	+	.	ID=mRNA00001;Parent=gene00001;Name=EDEN.1
chr1	.	mRNA	1050	9000	.	+	.	ID=mRNA00002;Parent=gene00001;Name=EDEN.2
chr1	.	mRNA	1300	9000	.	+	.	ID=mRNA00003;Parent=gene00001;Name=EDEN.3
chr1	.	exon	1300	1500	.	+	.	ID=exon00001;Parent=mRNA00003
chr1	.	exon	1050	1500	.	+	.	ID=exon00002;Parent=mRNA00001,mRNA00002
chr1	.	exon	3000	3902	.	+	.	ID=exon00003;Parent=mRNA00001,mRNA00003
chr1	.	exon	5000	5500	.	+	.	ID=exon00004;Parent=mRNA00001,mRNACHR1 mRNA00003
chr1	.	exon	7000	9000	.	+	.	ID=exon00005;Parent=mRNA00001,mRNA00002,mRNA00003

Genome Variation Format

```
##gvf-version 1.06
##genome-build GRCh37.1
##sequence-region chr16 1 88827254
chr16 UG SNV 291141 291141 33 + . ID=ID_1;Variant_seq=A,G;Reference_seq=G;
chr16 UG SNV 291360 291360 17 + . ID=ID_2;Variant_seq=G;Reference_seq=C;
chr16 UG SNV 302125 302125 67 + . ID=ID_3;Variant_seq=T,C;Reference_seq=C;
chr16 UG SNV 302365 302365 43 + . ID=ID_4;Variant_seq=G,C;Reference_seq=C;
chr16 UG SNV 302700 302700 75 + . ID=ID_5;Variant_seq=T;Reference_seq=C;
chr16 UG SNV 303084 303084 16 + . ID=ID_6;Variant_seq=G,T;Reference_seq=T;
chr16 UG SNV 303156 303156 90 + . ID=ID_7;Variant_seq=T,C;Reference_seq=C;
chr16 UG SNV 303427 303427 52 + . ID=ID_8;Variant_seq=T,C;Reference_seq=C;
chr16 UG SNV 303596 303596 66 + . ID=ID_9;Variant_seq=T,C;Reference_seq=C;
```

Sqid	Source	Type	Start	End	Score	Strand	Phase	Attributes
------	--------	------	-------	-----	-------	--------	-------	------------

Variant Annotation Tool (VAT)

- Adds functional annotation of the effect of sequence alterations (SNV, insertion, deletion) on sequence features (genes, mRNA)
- Takes as input a reference sequence (Fasta file), and set of gene models (GFF3 file) and a set of variants (GVF)
- Produces an annotated GVF file as output

```
VAT -a genome.fasta -f gene.gff3 variants.gvf > variants.vat.gvf
```

Variant_effect attribute

- Describes the effect of a sequence alteration on a sequence feature
 - The sequence_variant (the effect)
 - The Variant_seq allele index (which allele causes this effect)
 - The sequence_feature (what type of feature is affected)
 - The feature IDs (which features are affected)

```
Variant_seq=A,T;  
Variant_effect=missense_variant 0 mRNA NM_001160184 NM_032129;
```

Variant Selection Tool (VST)

- Applies complex set operations (intersection, union etc) to GVF files and produces a condensed representation of the genotypes.
- Takes as input a description of the set operation and a group of GVF files.
- Outputs population genotypes in CDR format.

```
VST -o 'I(0,1)' exome1.gvf exome2.gvf > affected.cdr
```

```
VST -o 'C(0,U(1,2))' kid.gvf mom.gvf dad.gvf > denovo.cdr
```

VST set operations

- **(U)nion:** All variants in all files.
- **(I)ntersection:** Variants shared by all files.
- **(C)omplement:** The left relative complement or variants unique to the first file (set).
- **(D)ifference:** The symmetric difference or variants unique to any one file (set).
- **(S)hared:** Variants shared by n files. $S(n,0..2);$
`'S(">2",0..2)').`
 - = Exactly n files share the variant.
 - > Greater than n files share the variant.
 - < Less than n files share the variant.

VAAST Condenser File (CDR)

Seqid	Start	End	Type	Effect	Reference	Genotypes	
chr1	877831	877831	SNV	missense_variant	T W	0-3 C:C R:R	
chr1	881627	881627	SNV	synonymous_variant	G L	0,2-3 A:A L:L	1 A:G L:L
chr1	881918	881918	SNV	missense_variant	G S	2 A:G L:S	
chr1	887801	887801	SNV	synonymous_variant	A T	1 A:G T:T	0,2-3 G:G T:T
chr1	888639	888639	SNV	synonymous_variant	T E	0 C:C E:E	1 C:T E:E
chr1	888659	888659	SNV	missense_variant	T I	0,1-2 C:C V:V	3 C:T V:I
chr1	889238	889238	SNV	missense_variant	G A	1 A:G V:A	
chr1	897325	897325	SNV	synonymous_variant	G A	0,2 C:C A:A	1 C:G A:A
chr1	897738	897738	SNV	synonymous_variant	C L	1 C:T L:L	
chr1	900505	900505	SNV	synonymous_variant	G V	3 C:C V:V	2 C:G V:V
chr1	900972	900972	SNV	3_prime_UTR_variant	T	0,2 G:G	1 G:T
chr1	901023	901023	SNV	3_prime_UTR_variant	T	0,2-3 C:C	1 C:T
##	GENOME-LENGTH	914121104					
##	GENOME-COUNT	7					
##	GENDER	F:0-1,3	M:2				
##	FILE-INDEX	0		A12.vat.gvf			
##	FILE-INDEX	1		B34.vat.gvf			
##	FILE-INDEX	2		C56.vat.gvf			
##	FILE-INDEX	2		D78.vat.gvf			

VAAST

- Scores and prioritizes features in a probabilistic fashion for their likelihood of being associated with a disease phenotype.
- Takes as input a set of gene models (GFF3), a set of variants for background/healthy genomes (CDR) and a set of variants for target/disease genomes (CDR).

```
VAAST -m lrt -o Output_name genes.gff3 background.cdr target.cdr
```

VAAST Uses Variant Frequencies in a Probabilistic Fashion

Composite Likelihood Ratio Test

$$\lambda = \ln \left(\frac{L_{Null}}{L_{Alt}} \right)$$

Maximum Likelihood
of the Null Model
(No Difference)

Maximum Likelihood
of the Alternate Model
(There is Difference)

VAAST Uses Variant Frequencies in a Probabilistic Fashion

$$\lambda = \sum_{i=1}^k \ln \left(\frac{n_i \hat{p}_{Yi}^{B_{Yi} + T_{Yi}} (1 - \hat{p}_{Yi})^{B_{Xi} + T_{Xi}}}{a_i \hat{p}_{BYi}^{B_{Yi}} (1 - \hat{p}_{BYi})^{B_{Xi}} \hat{p}_{TYi}^{T_{Yi}} (1 - \hat{p}_{TYi})^{T_{Xi}}} \right)$$

Annotations:

- LRT for AAS (HGMD/1KG) points to $n_i \hat{p}_{Yi}^{B_{Yi} + T_{Yi}} (1 - \hat{p}_{Yi})^{B_{Xi} + T_{Xi}}$
- ML of Null points to the denominator of the fraction.
- ML of Alt points to the numerator of the fraction.

- p: MAF
- B/T: Background/Target Allele Counts
- X/Y: Minor/Major Allele

Conservation-controlled amino acid scoring matrix – CASM

- The amino acid severity parameter is adjusted to account for the degree of phylogenetic conservation at a site.
- Conservation data comes from PhastCons scores which estimate a probability of a site being under negative selection (conserved).
- Each AAS type is scored at all sites with a PhastCons score of 0 and 1 (the two extremes).
- The AAS severity parameter at any PhastCons score is then linearly interpolated between those extremes.
- The effect is that the AAS severity parameter diminishes with diminishing conservation.

VAAST Uses Variant Frequencies in a Probabilistic Fashion

- VAAST gives us the likelihood of the composite genotype of a given gene in the target given the background.
- Do allele and AAS frequencies differ between background and target genomes within a given gene or feature?
- Composite likelihood calculation assumes independence across sites. To control for LD, statistical significance is estimated by permutation test.
- Multiple test correction for number of features (~20,000) is two orders of magnitude better than for the number of variants (~3,500,000).

VAAST: highly accurate variant prioritization

	Percent Judged Deleterious					
	Dream tool	SIFT ²	ANNOVAR ³	PolyPhen2	Mutation Taster	VAAST
Disease alleles ^A	100%	58%	71%	84%	84%	99%
Healthy alleles ^B	0%	12%	1%	16%	16%	10%

Accuracy (Sn + Sp)/2	100%	80%	88%	86%	86%	95%
-------------------------	------	-----	-----	-----	-----	-----

^A1454 high quality, published disease-causing/predisposing OMIM alleles¹

^B1454 Variants randomly selected from 5 different healthy CEU individuals' genomes

VAAST Filters

■ Inheritance model

- Dominant: Only score one allele per feature
- Recessive: Only score two alleles per feature

■ Locus heterogeneity

- Require all affected individuals to have a scoring allele

■ Complete penetrance

- Don't score a feature if anyone in the background shares it's scoring alleles.

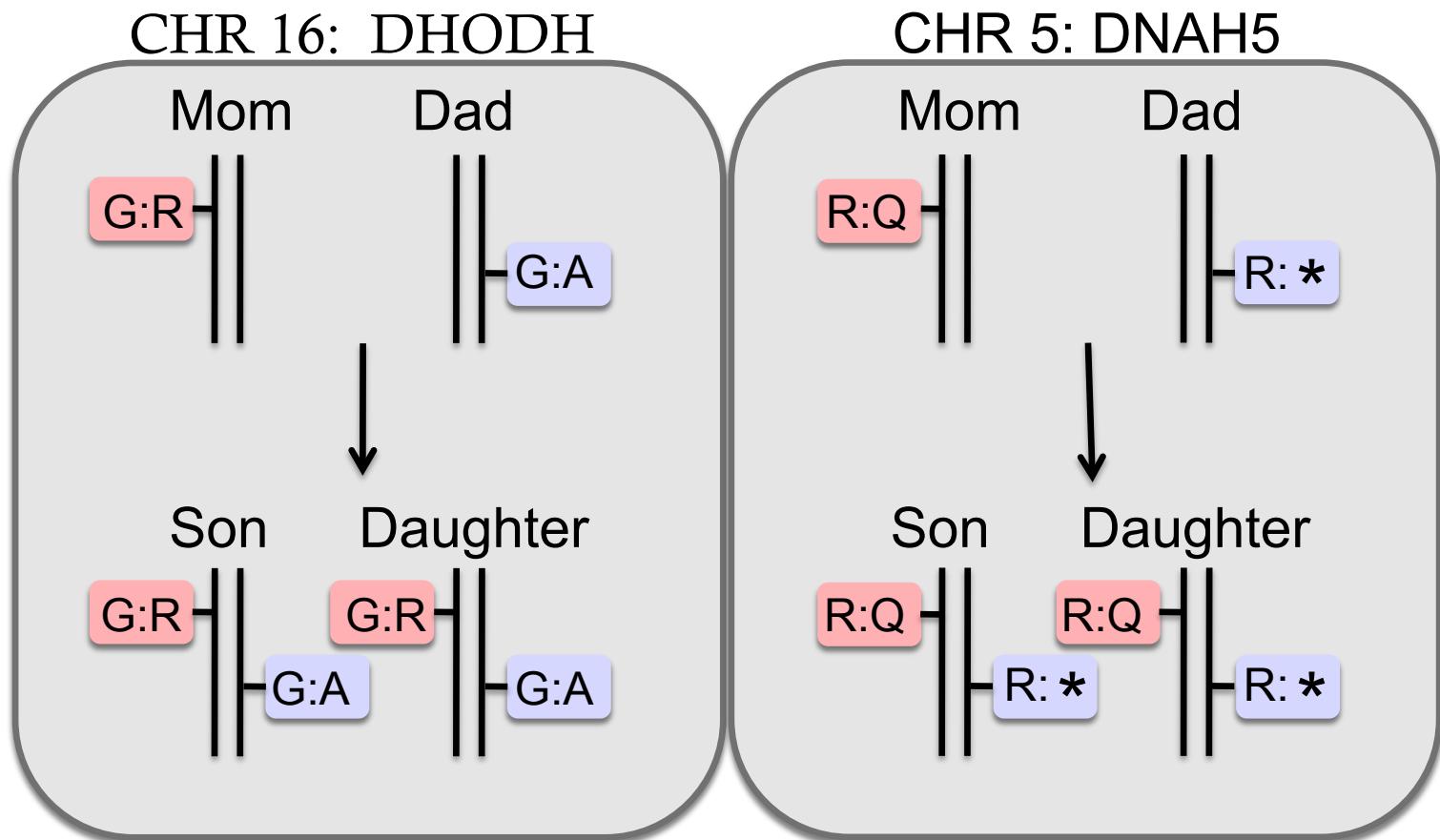
■ Rate/PAR

- Hold the MAF in the background/target below the given value.

Dangers of filtering

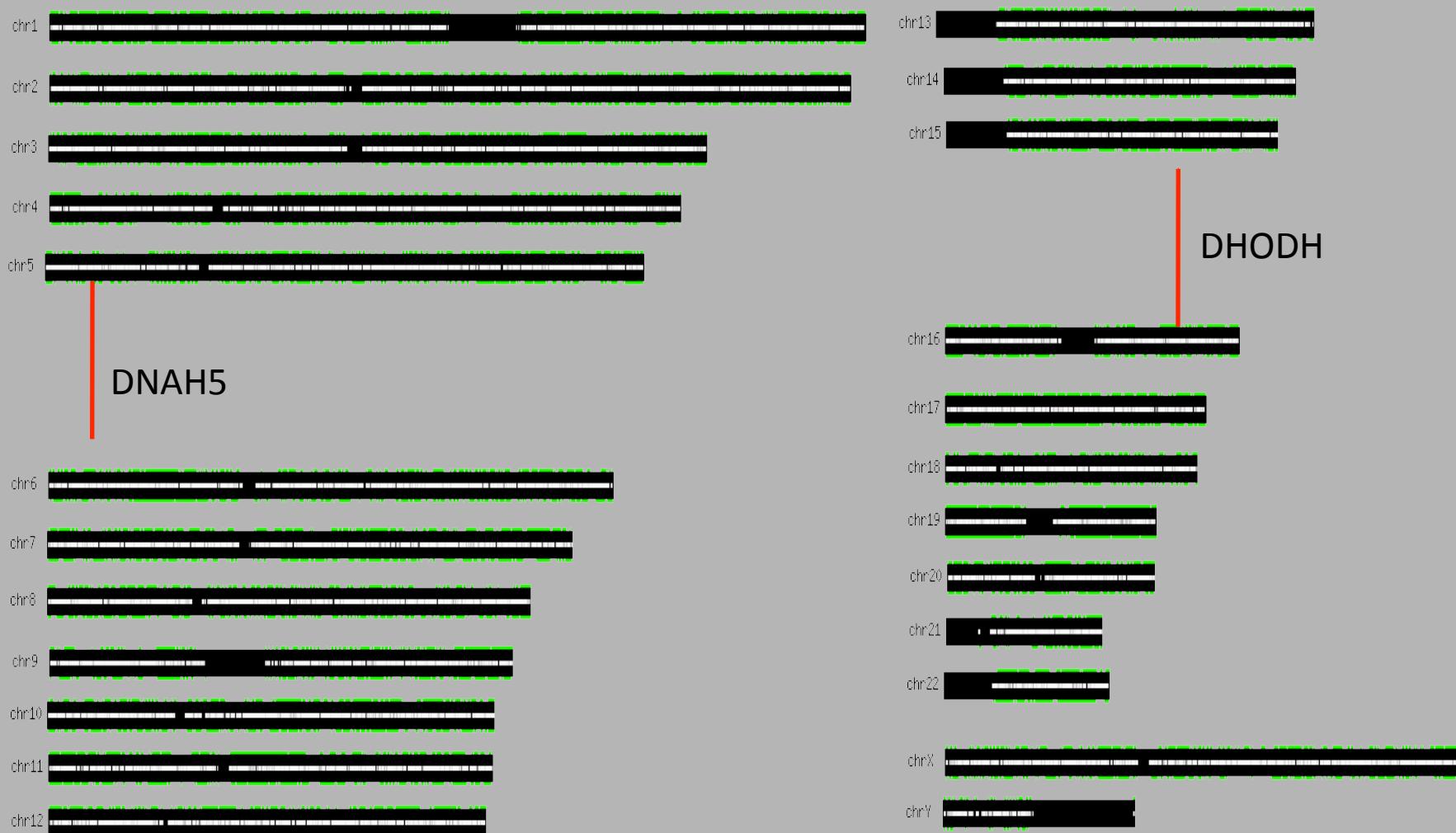
- Filters are binary – you either pass or fail.
- Filters (usually) don't consider missing data.
- Filters aren't able to incorporate multiple factors into the equation.

Alleles Responsible for Miller Syndrome in Utah Kindred



• Ng *et al*, Nature Genetics 42, 30–35, 2010
• Roach, *et al*, Science 328 636, 2010

Schematic of VAAST Analysis of Utah Miller Kindred Using a Single Quartet

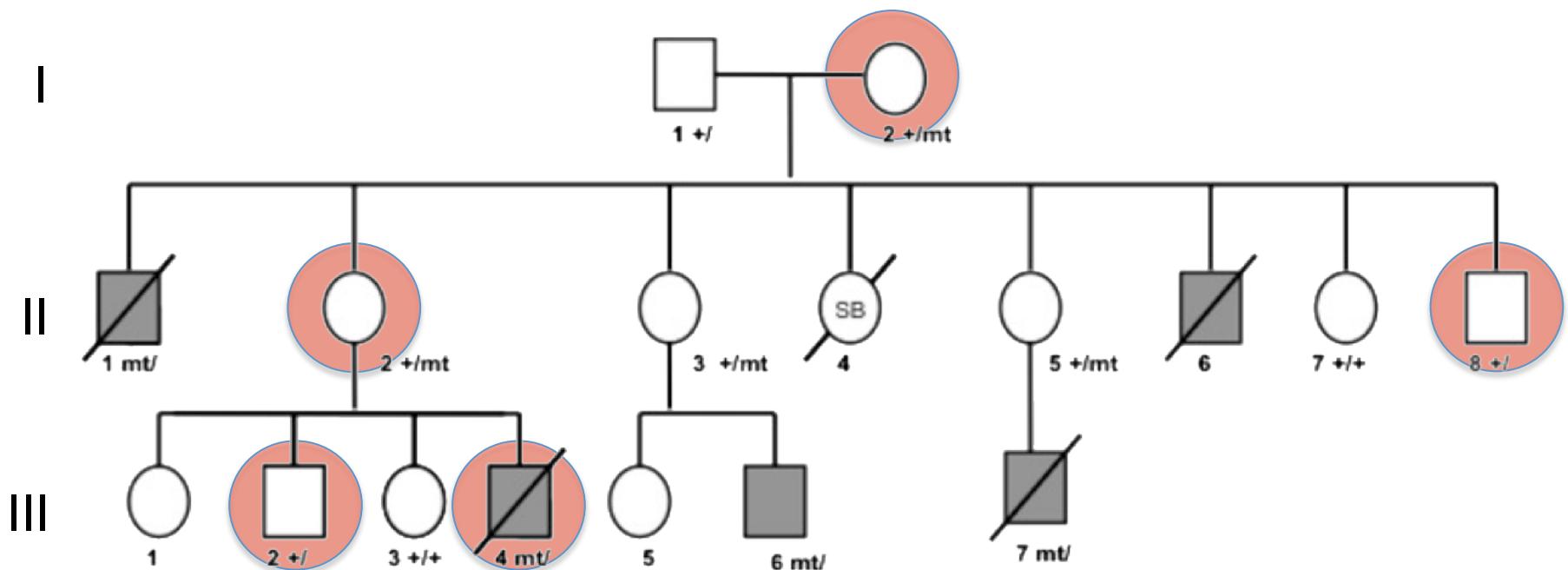


A rare X-linked mendelian disorder

- A Utah family coming to the University Hospital for 20+ years
- About half of the male offspring die around 1 year of age
- Aged appearance
- Craniofacial anomalies
- Hypotonia
- Global developmental delays
- Cardiac arrhythmias

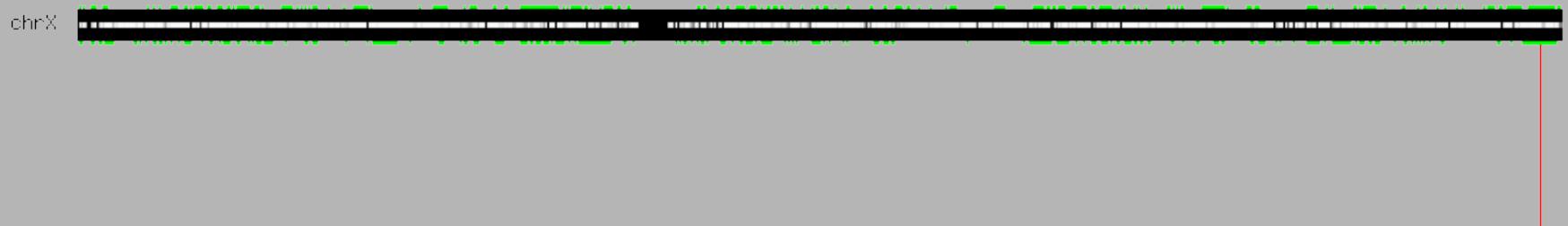


Four Affected Boys over Two Generations



Identifying Candidate Genes

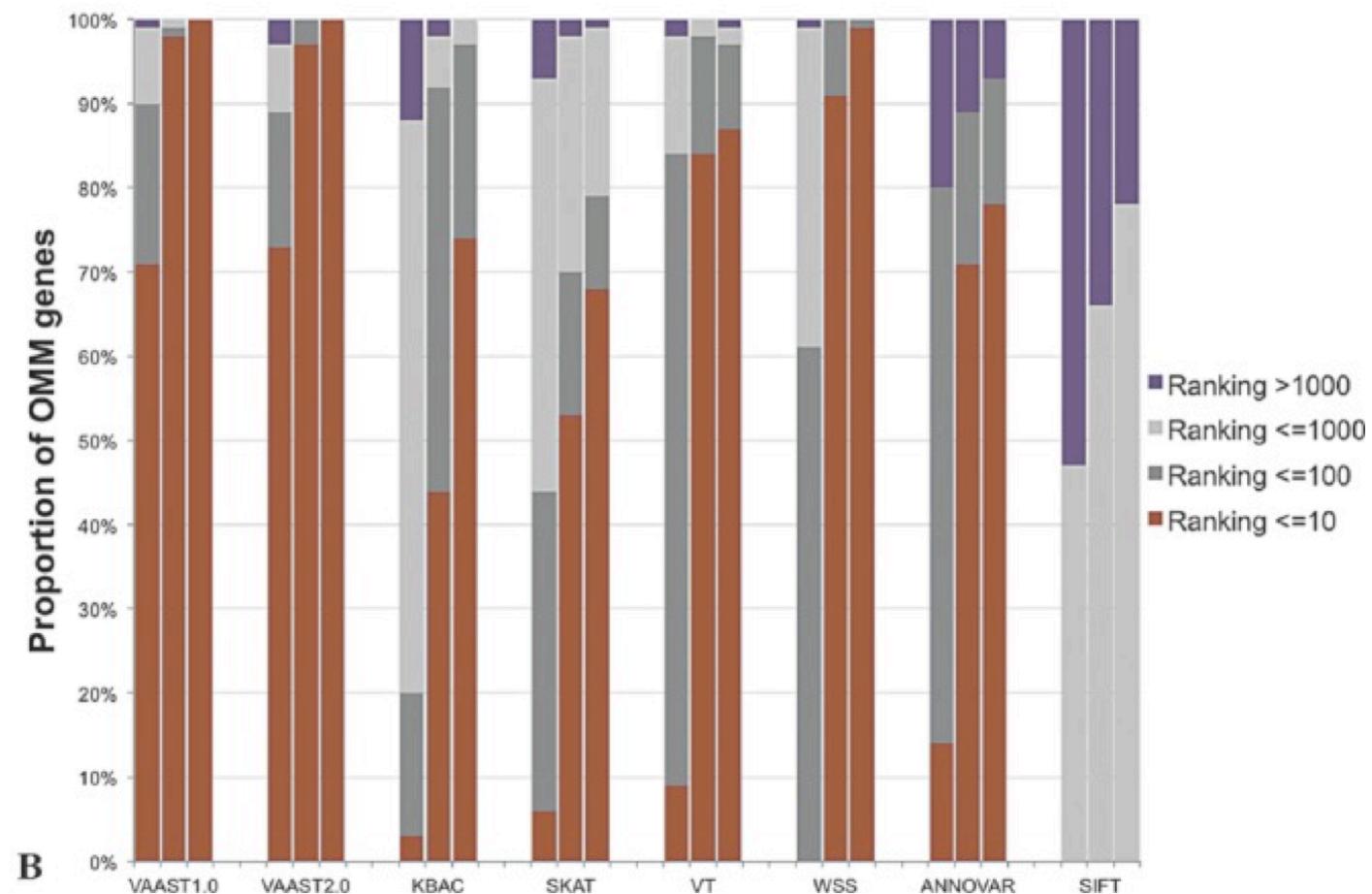
- VAAST identified NAA10 as candidate gene
 - Run entire pipeline in an afternoon
 - 3 candidate genes (NAA10 ranked 2) proband only
 - 1 candidate gene (NAA10) with pedigree



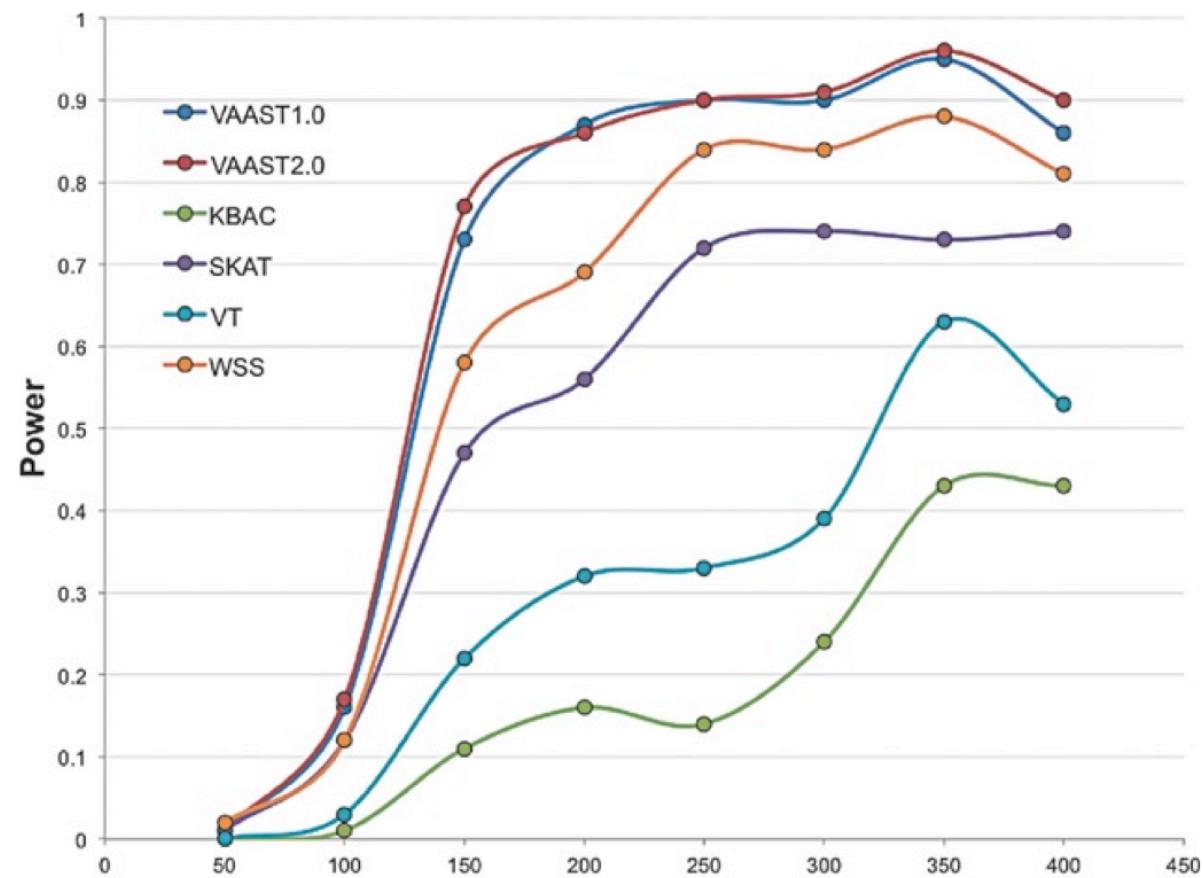
VAAST benchmark – 100 OMIM diseases

- Randomly choose known disease genes from OMIM
- Randomly insert one or more published disease causing variants for that gene into a personal exome
- Assay the ability of VAAST, SIFT1 and ANNOVAR2 to identify the disease gene in a genome-wide screen
- Repeat for 100 different genes under a variety of different scenarios, e.g. dominant, recessive, various case cohort sizes etc.

VAAST benchmark – 100 OMIM Diseases



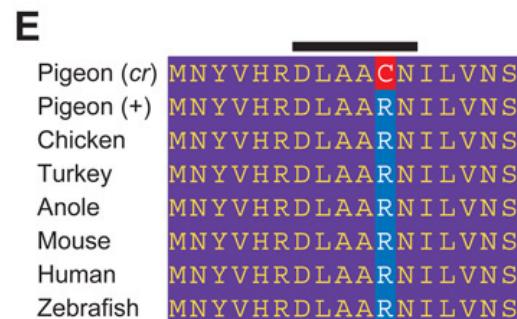
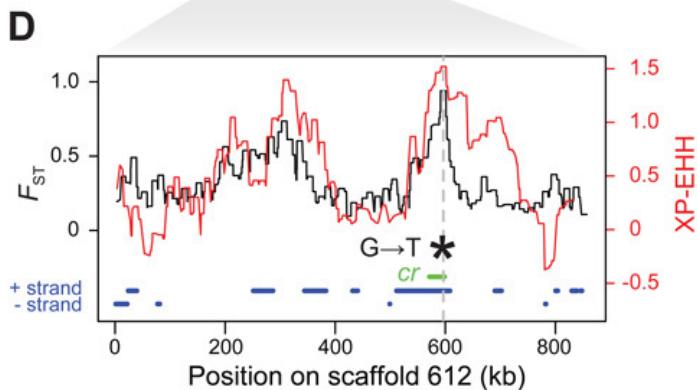
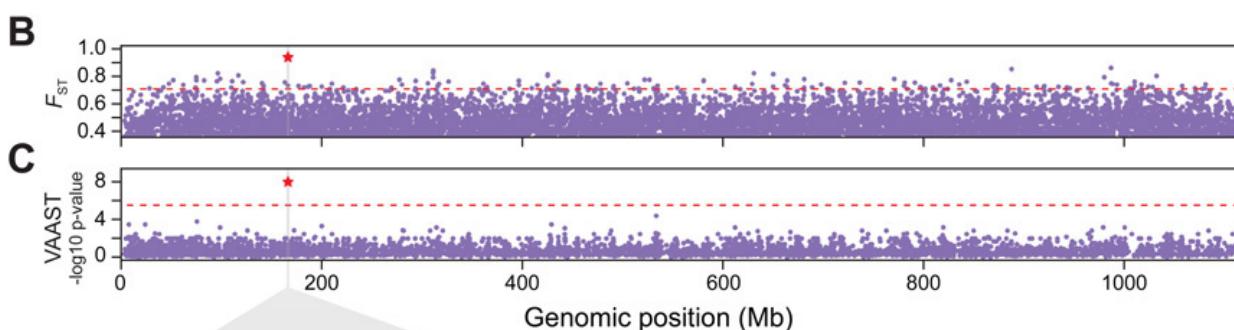
VAAST benchmark – common disease



Power comparisons over published LPL
(Power is calculated based on 100 bootstraps)

Hu et al. Genet Epidemiol. 2013

VAAST in non-human organisms

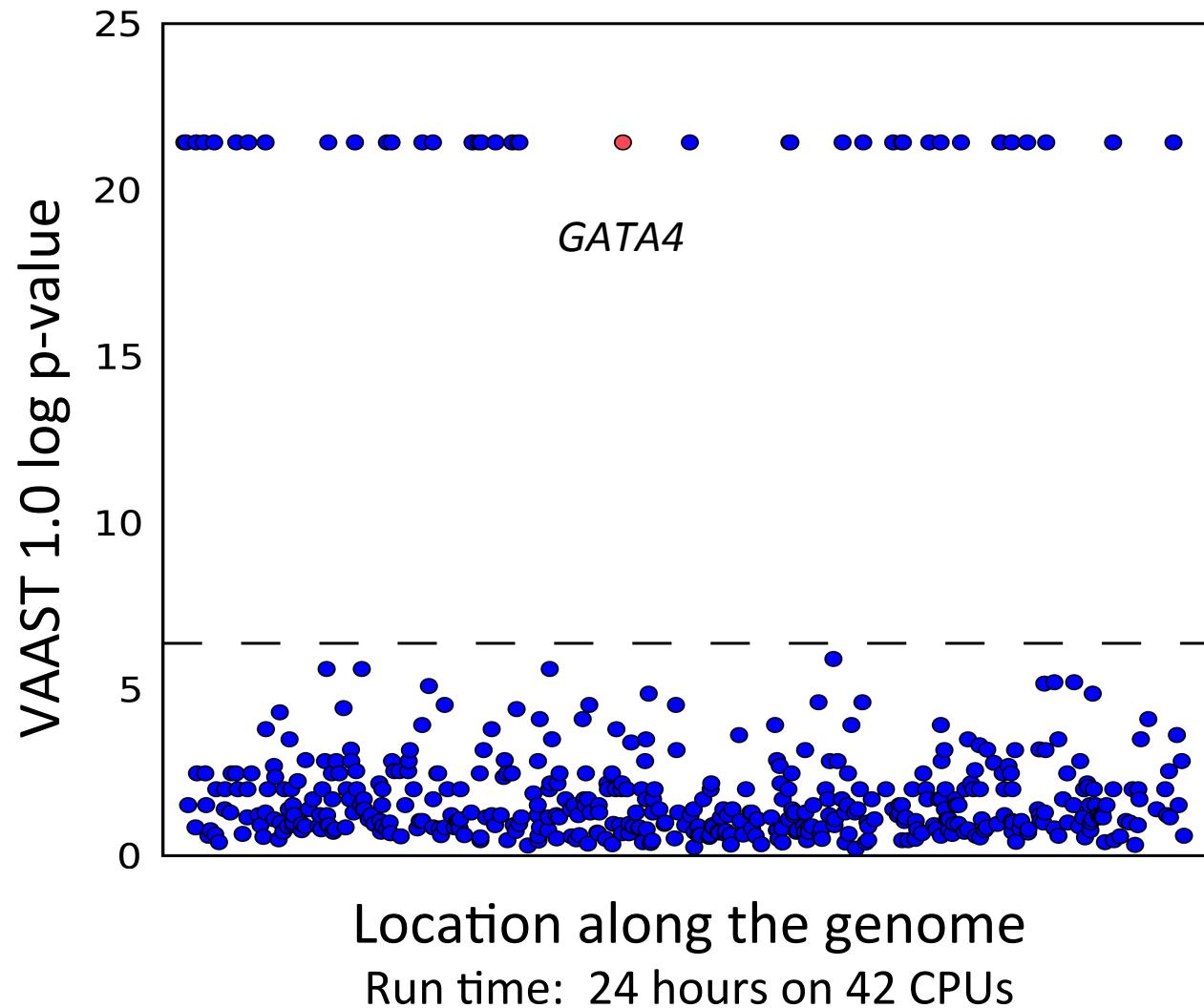


Shapiro et al. Science. 2013

pVAAST for pedigree analyses

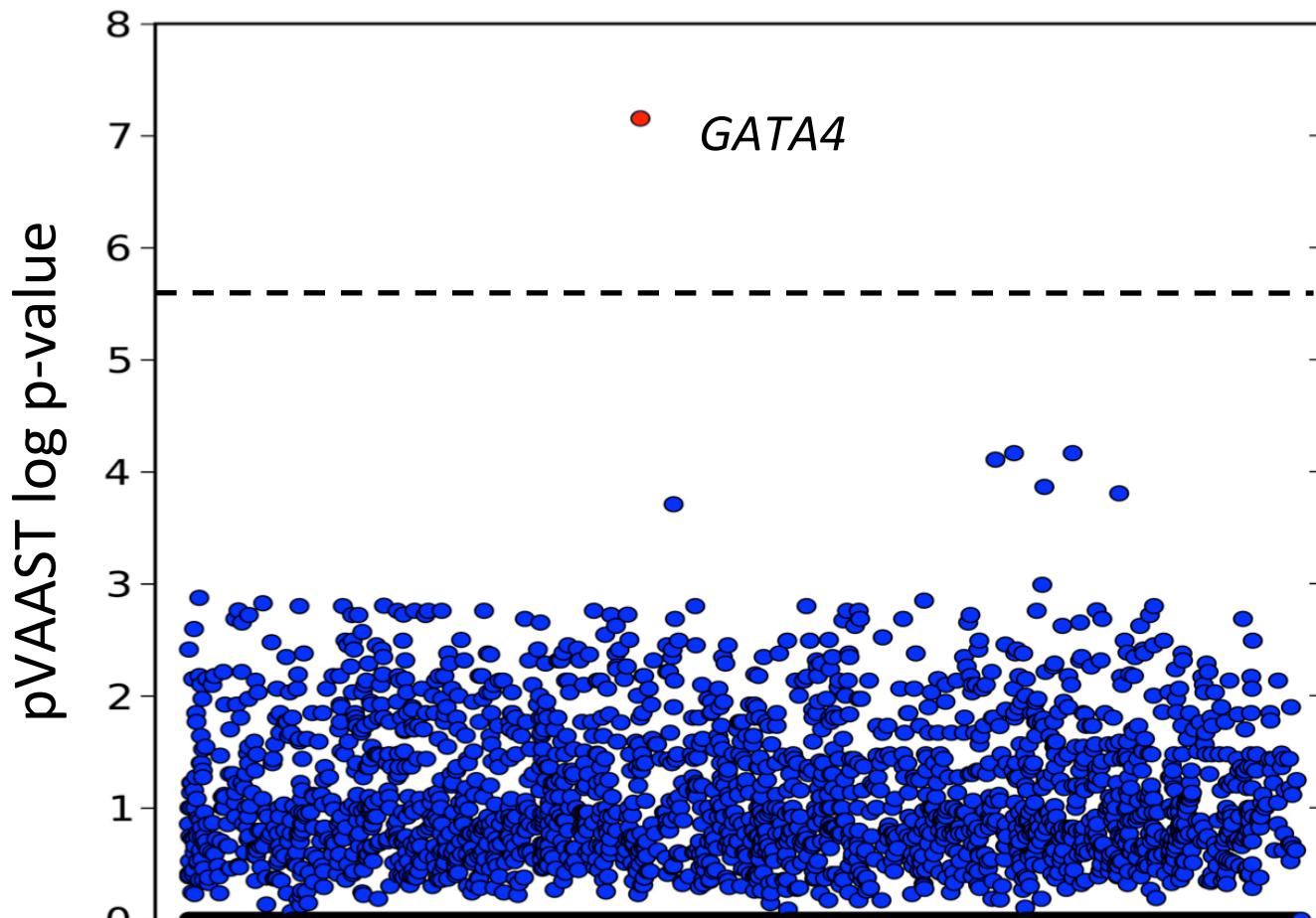
- Extends VAAST to incorporate family data (pedigrees)
- pVAAST performs linkage analysis by calculating a gene-based LOD designed for NGS
- The LOD score at each locus is incorporated directly into the CLRT.
- In large-scale simulation studies and re-analysis of known disease pedigrees, pVAAST had significantly higher statistical power compared other tools – including VAAST.

VAAST 1.0 cardiac septal defect



Hu et al submitted, Original study: Garg et al. Nature. 2003

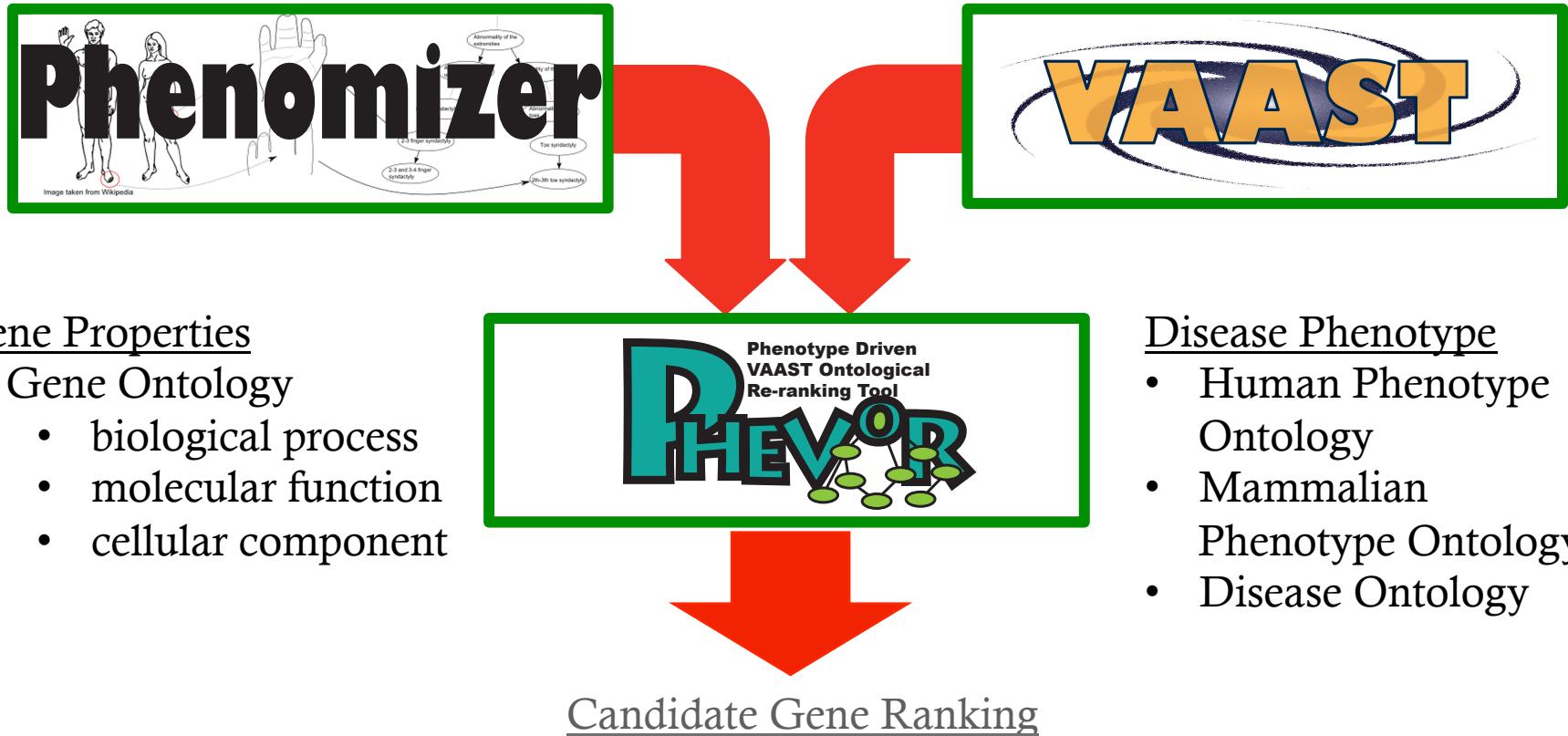
pVAAST - cardiac septal defect



Location along the genome

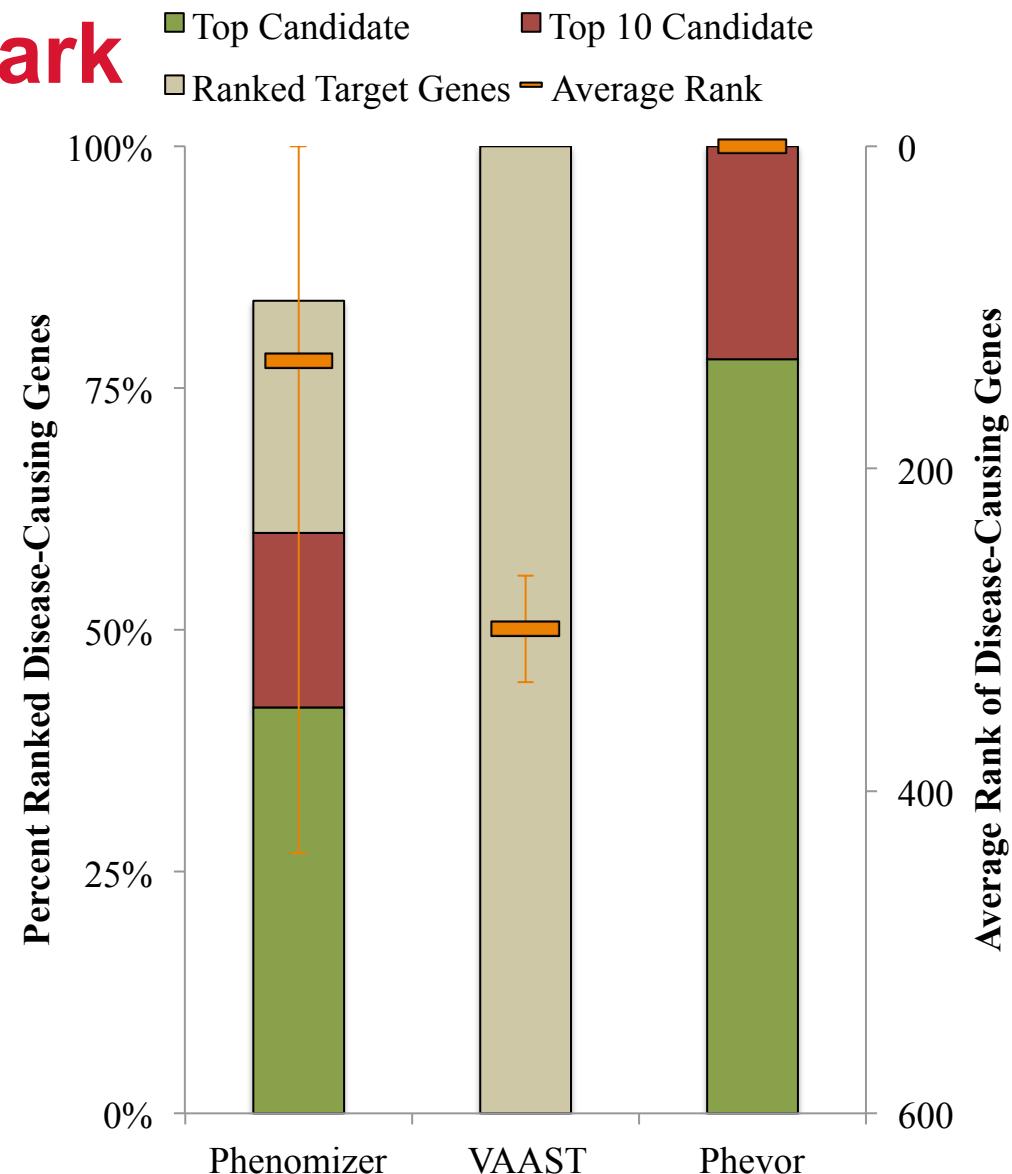
Run time: 4 hours on 42 CPUs

Phevor



Phevor benchmark

- 50 Dominant Disease-Causing Variants
- 50 Corresponding Phenomizer Reports
- Spiked into a single healthy exome
- Average Ranks
 - Phenomizer – 134*
 - VAAST – 300
 - Phevor – 1.3



Omicia - Opal

Omicia Opal - 0.8.0 (beta)

VAAST Quad Report

Overview
 Affected Child: 201139
 Affected Sibling: 201141
 Unaffected Mother: 201138
 Unaffected Father: 201140
 Background: 1K Project Phase 1
 VAAST Release: RC1.0 (recessive model)

Note	Hide	Class	Gene	Position dbSNP	Change	Proband Zygosity	Sibling Zygosity	Father Zygosity	Mother Zygosity	Effect	Global MAF	Omicia Score	V-Score	G-Score	Evidence
Cat 1	DNAH5	chr5 13864742	c.4360C>T p.Arg1454*	het	het	-	het	stop gained	G100% A:0%	0.817	28.8	54.01	upstream from hgmd: Primary ciliary dyskinesia (pubmed, omim)		
Cat 3	DNAH5	chr5 13792155	c.8396G>A p.Arg2799Gln	het	het	het	-	non-synon	-	0.875	26.02	54.01			
Cat 1	DHODH	chr16 72050942	c.454G>A p.Gly152Arg	het	het	-	het								
Cat 1	DHODH	chr16 72055110	c.605G>C p.Gly202Ala	het	het	-	het								
Cat 3	KIAA0556	chr16 27784497 rs117316067	c.4276G>A p.Glu1426Lys	het	het	-	het								
Cat 3	KIAA0556	chr16 27784497 rs117316067	c.4685G>A	het	het	-	het								

OMICIA Opal - Dev 0.10.0

VAAST Trio Report

Overview
 Proband: 200571
 Unaffected Mother: 200572
 Unaffected Father: 200573
 Background: 1K Project Phase 1
 VAAST Release: RC1.0 (recessive model)

Note	Hide	Class	Gene	Position dbSNP	Change	Proband Zygosity	Father Zygosity	Mother Zygosity	Effect	Global MAF	Omicia Score	V-Score	G-Score	Evidence
3 - VUS	KRT24	chr17 38858135 rs11309872	c.666_668delT p.Asn222fs	hom	het	het	het	frameshift deletion	-	0.408	23.54	23.54		
4 - LB	CR1	chr1 207790088 rs3811381	c.1348C>G p.Leu450Val	hom	het	-				0.087	23.53	23.53	hgmd-dp: Idiopathic pulmonary fibrosis, assoc. with ? (pubmed, omim)	
1 - P	DHODH	chr16 72055110	c.605G>C p.Gly202Ala	het	-	het	het			0.87	9.56	20.26	omim: Miller Syndrome (omim)	
1 - P	DHODH	chr16 72050942	c.454G>A p.Gly152Arg	het	het	het	het			0.938	12.7	20.26	hgmd-dm: Miller syndrome (pubmed, omim)	
4 - LB	CYP2A7	chr19 41384781 rs112946838	c.715A>C p.Lys239Gln	het	-					0.277	10.17	17.61	omim: Miller Syndrome (omim)	
4 - LB	CYP2A7	chr19 41383153 rs261144	c.1103T>C p.Met368Thr	het	het					0.297	9.44	17.61	hgmd-dm: Miller syndrome (pubmed, omim)	

Filter VAAST Results
 Exclude selected Polymorphic Set:
 All

Show only selected Gene Set:
 VAAST Disease Genes

Filter By VAAST Gene Score:
 15 100 G:3176

Personal Variants in this Gene

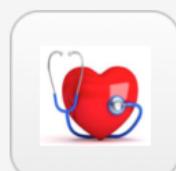
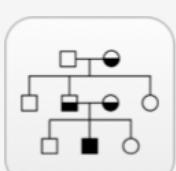
Position	Transcript	Transcript HGVS	Protein	Pro
72042682	NM_001361.3	c.194>C	NP_001352	p.Ly
72050942	NM_001361.3	c.454G>A	NP_001352	p.G
72055110	NM_001361.3	c.605G>C	NP_001352	p.G

Page 1 of 1 Displaying 1 to 6 of 6 items

Omicia Home - Terms of Service - Privacy Policy - Blog
 © 2012, Omicia, Inc. All rights reserved.

Development supported by NIH SBIR grants 1R4HG003667 to Omicia/Yandell,
 SBIR 1R44HG002991 to Omicia

VAAST in Opal

	Omicia Variant Miner ***** 225 Ratings		Omicia/University of Utah VAAST Solo Analysis ***** 67 Ratings		Omicia/University of Utah VAAST Cohort Analysis ***** 11 Ratings
	Omicia Flex Trio Analysis ***** 132 Ratings		Omicia/University of Utah VAAST Duo Analysis ***** 44 Ratings		Omicia Gene Health Analysis ***** 14 Ratings
	Omicia Flex Quad Analysis ***** 55 Ratings		Omicia/University of Utah VAAST Trio Analysis ***** 17 Ratings		Omicia Variant Load Analysis ***** 19 Ratings
	Omicia Flex Complex Analysis ***** 12 Ratings		Omicia/University of Utah VAAST Quad Analysis ***** 7 Ratings		Omicia Panel Reporter ***** 6 Ratings

VAAST in Summary

- Probabilistic Disease Gene Finder
- Feature Based
- Both Allele and AAS Frequencies
- Considers the Inheritance Model
- As few as 1-2 target genomes can be sufficient to identify causative gene.
- Complete analysis pipeline
- Many parameters allow fine-grained control of analysis
- pVAAST and Phevor on the way for otherwise under-powered analyses.

Yandell Lab – Variant Annotation, Analysis and Search Tool

Yandell Lab – Variant Annotatio... +

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

illumina sequenc... Bookmarks

Most Visited Barry's Wiki Yandell Lab Home Yandbeck Wiki The Sequence O...

Yandell Lab

Department of Human Genetics - University of Utah

Home People Research Software MWAS Publications About Links Utah Contact Internal

Variant Annotation, Analysis and Search Tool

VAAST (the Variant Annotation, Analysis and Search Tool) is a probabilistic search tool for identifying damaged genes and their disease-causing variants in personal genome sequences. VAAST builds upon existing amino acid substitution (AAS) and aggregative approaches to variant prioritization, combining elements of both into a single unified likelihood-framework that allows users to identify damaged genes and deleterious variants with greater accuracy, and in an easy-to-use fashion. VAAST can score both coding and non-coding variants, evaluating the cumulative impact of both types of variants simultaneously. VAAST can identify rare variants causing rare genetic diseases, and it can also use both rare and common variants to identify genes responsible for common diseases. VAAST thus has a much greater scope of use than any existing methodology.

Publications

[A probabilistic disease-gene finder for personal genomes](#)
Yandell M Huff CD Hu H Singleton M Moore B Xing J Jorde L Reese MG
Genome Res. 2011 Jul

[Using VAAST to Identify an X-Linked Disorder Resulting in Lethality in Male Infants Due to N-Terminal Acetyltransferase Deficiency](#)
Rope AF Wang K Ejventh R Xing J Johnston JJ Swensen JJ Johnson WJ Moore B Huff CD Bird LM Carey JC Opitz JM Stevens CA Jiang T Schank C Fain HD Robison R Dalley B Chin S South ST Pysher TJ Jorde LB Hakonarson H Lillehaug JR Biesecker LG Yandell M Arnesen T Lyon GJ
Am J Hum Genet. 2011 Jul 15;89(1):28-43

Acknowledgements

VAAST Development

- Chad Huff
- Hao Hu*
- Lynn Jorde*
- Edward Kirulata*
- Marco Falcioni*
- Barry Moore*
- Martin Reese
- Jinchuan Xing
- Mark Yandell

Sequence Ontology

- Mike Bada
- Colin Batchelor
- Karen Eilbeck
- Barry Moore
- Shawn Reyneerson

Yandell Lab

- Michael Campbell
- Daniel Ence
- Steven Flygare
- Carson Holt
- Brett Kennedy
- Zev Kronenberg
- Qing Li
- Gordon Lemmon
- Barry Moore
- EJ Osbourne
- Scott Watkins
- Mark Yandell

Omicia

- Marco Falcioni
- Edward Kirulata
- Martin Reese
- Jeff Rule
- Charlene Rigby
- Andrew Gao

Funding

- NHGRI

Acknowledgements



Base Quality Score

- Base-calling considers many possible sources of error in the sequencing process:
 - Mixed clusters
 - Out of phase clusters
 - Overlapping emission spectra
- The base quality score (BQS) is based on the probability that the base was called wrong.
- The quality score is "Phred scaled"
 - $-10 * \log_{10}(\text{probability of miscall})$

$$-10 * \log_{10}(0.0013)) + 33 = 72 \Rightarrow H \text{ (ASCII)}$$

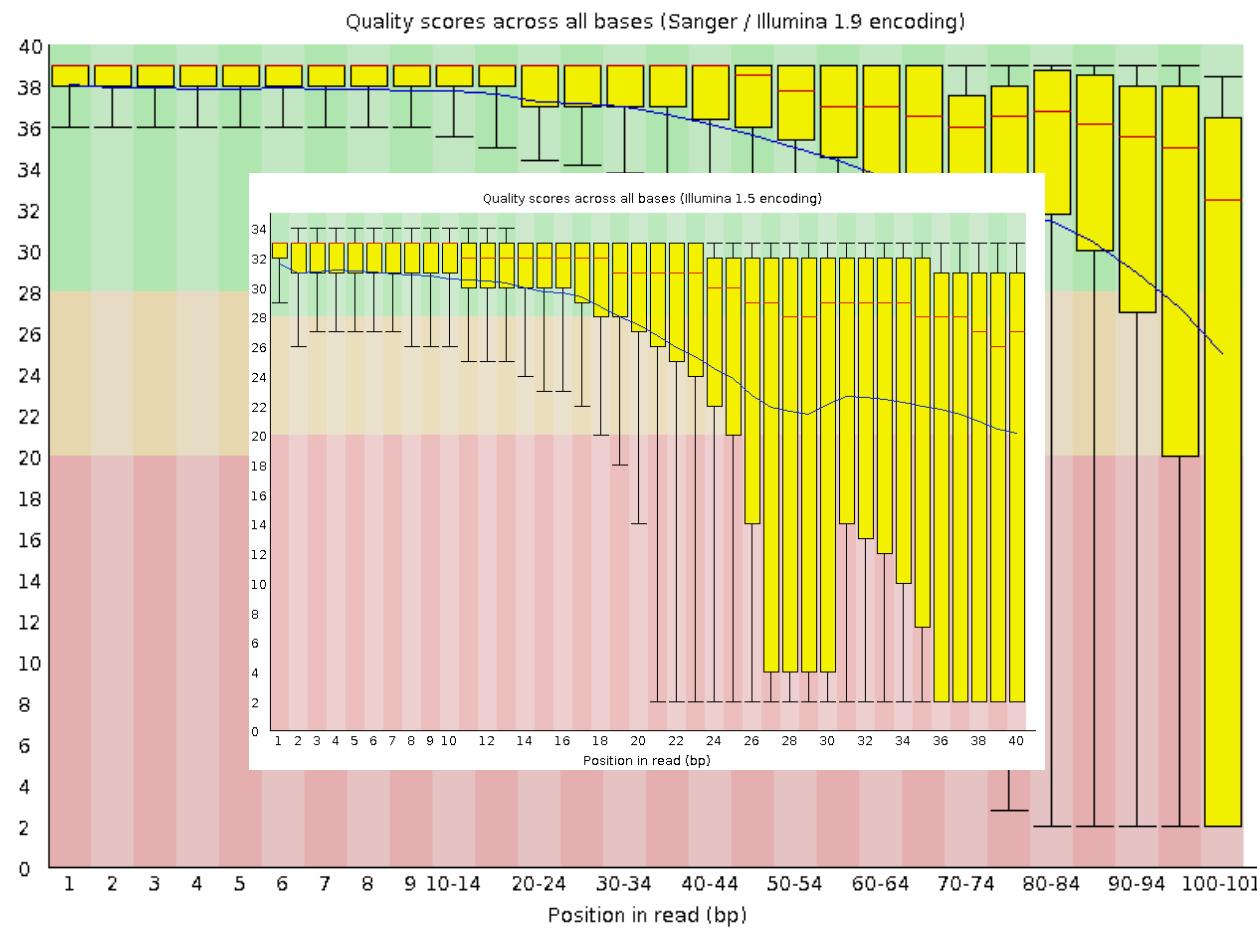
PHRED Scaled

FastQ Encoded

FastQ

```
@HWI-ST179R:404:D1YNUACXX:7:1101:1228:1937 1:N:0:ATCACG
NTGATACAGCTGAACTGAACTCCTGCCTGTACAGCTCGTTCTACAAGATTCCAGACCTGGAAGATGATGGC
+
#1=DDDFFGHHHJJHHIIHIIIFHJJFHCHGIJJGIICGIJIGIJJEHIAGHIIGJJGIIJIHI:DGCHIG
@HWI-ST179R:404:D1YNUACXX:7:1101:1184:1940 1:N:0:ATCACG
NATGTCAGCCCACCCAGGAGAACACAGACCCAAGGGAACCCCCACTCCAGCCAAGAGAAGCCGTGAGTGAA
+
#1=BADDHHHFIIIIIIII>GIHIIIIIIIIIEIGGHIIIFHHEEECCCCB9>=?>CD
@HWI-ST179R:404:D1YNUACXX:7:1101:1193:1964 1:N:0:ATCACG
NTGCCACAGGGCGGTGTAAGACAGGAGTCCATCTGGGCAGGGTGAGAGGATGGGGTCAGAGGCACTAA
+
#1=DDFFFHHDHFJJ6@<FHJGIJJIIIC>EEEHHEFFFDDDD(8?BADD??@BDD07ACDDDB@B<CC
@HWI-ST179R:404:D1YNUACXX:7:1101:1166:1977 1:N:0:ATCACG
CTTCTTGCACCTCAAGGGATCACTCCCTCTAGGCCGTTGCCATTCTCGCTGGAAACCCTCT
+
@?<D:DDDDHHFHGIIIF;DGGAHIIIGHHHIIIDHIIIDH9B?CCB=?@?B:::>:@3>>@C@CCB<@#
@HWI-ST179R:404:D1YNUACXX:7:1101:1423:1940 1:N:0:ATCACG
NAGATGTTGGCTATAGGAATACGGCAGGAAAATGAAACGTTGTGCATGGCAGGGCAGCATCACTTGGGGATC
+
#11ADDDDHHDHICFHGHIIHHIIII7DC?EH@B?C@6;<B;3=CC:AAC5>?B9?&8?<?#####
```

FastQC Report



Sequence Alignment

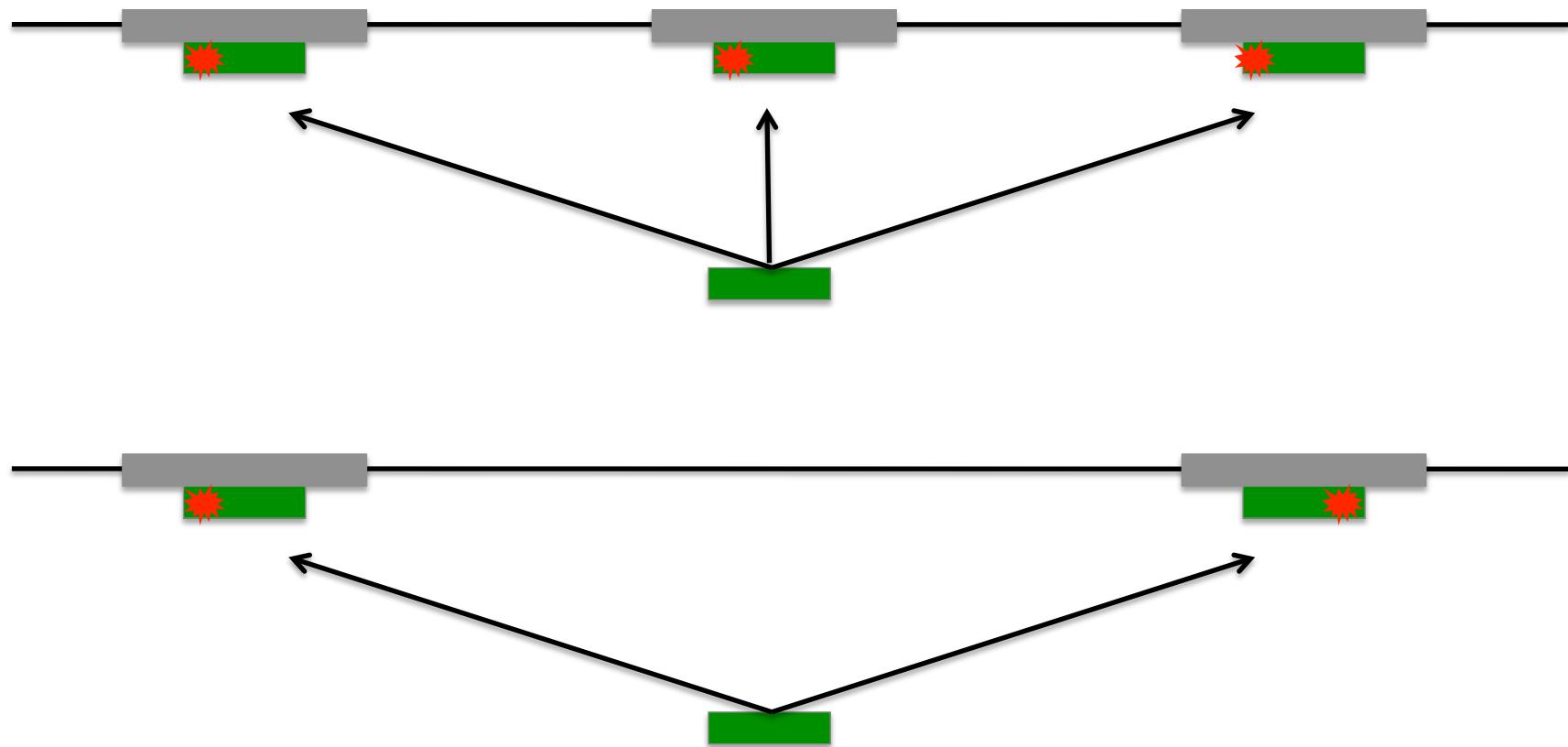
- BLAST
 - Fast (at least we used to think so)
 - Accurate – Smith-Waterman algorithm guarantees optimal alignment.
 - Seed and extend (SW dynamic)
- Short-read aligners
 - Fast
 - Sort of accurate – Many heuristics
 - Seeding and extension (heuristics terminate unlikely extensions).

SAM/BAM ish read)

Challenges for short-read aligners

- Repetative regions
 - Multiple alignments with the same score
 - How many of multiple alignments to score
 - Which SNV should I choose
- Insertions and deletions
 - Aligners hate to open gaps
 - Working one read doesn't allow context
 - Indel edges are inconsistent
 - A halo of bogus SNVs can be induced

Multi-mapping reads



Challenges for short-read aligners

- Repetative regions
 - Multiple alignments with the same score
 - How many of multiple alignments to score
 - Which SNV should I choose
- Insertions and deletions
 - Aligners hate to open gaps
 - Working one read at a time doesn't allow context
 - Indel placements are inconsistent
 - A halo of bogus SNVs can be induced

Indel alignment

Reference Sequence

TACAGCTGAACTGAACTCCTGCCTGTACAGCTCGTTTCTACAAGATTCCAGACCTGGAA

TGCCTGTACAGCTCGTT - TCTACAAGATT

AACTGAACTCCTGCCTGTACAGCTCGT

TGAAACTCCTGCCTGTACAGCTCG - TTTCTA

TGTACAGCTCGTT - CTACAAGATTCCAGA

CTCCTGCCTGTACAGCTCGTTTCTACAAG

ACTCCTGCCTGTACAGCTCGTT C - TACAA

GCCTGTACAGCTCGT - TTCTACAAGATTCC

Polishing Alignments

- Base quality recalibration – covariant analysis
 - Readgroup
 - Cycle
 - Sequence context
- Local re-alignment around indels (Local assembly)
 - Improves indel calls
 - Reduces SNV false positives

Covariant based BQS recalibration

Reference Sequence

TACAGCTGAACTGAACTCCTGCCTGTACAGCTCGTTTCTACAAGATTCCAGACCTGGAA

TGCCTGTACAGCTCGTTTCTACGAGATT

AACTGAACTCCTGCCTGTACAGCTCGT

TGAAACTCCTGCCTGTACAGCTCGTTTCTA

TGTACAGCTCGTTTCTACGAGATTCCAGA

CTCCTGCCTGTACAGCTCGTTTCTACGAG

ACTCCTGCCTGTACAGCTCGTTTCTACG

GCCTGTACAGCTCGTTTCTACGAGATTCC

Polishing Alignments

- Base quality recalibration – Covariant analysis
 - Readgroup
 - Cycle
 - Sequence context
- Local re-alignment around indels (Local assembly)
 - Improves indel calls
 - Reduces SNV false positives

Indel alignment

Reference Sequence

TACAGCTGAACTGAACTCCTGCCTGTACAGCTCGTTTCTACAAGATTCCAGACCTGGAA

TGCCTGTACAGCTCGTT - TCTACAAGATT

AACTGAACTCCTGCCTGTACAGCTCGT

TGAAACTCCTGCCTGTACAGCTCG - TTTCTA

TGTACAGCTCGTT - CTACAAGATTCCAGA

CTCCTGCCTGTACAGCTCGTTTCTACAAG

ACTCCTGCCTGTACAGCTCGTT C - TACAA

GCCTGTACAGCTCGT - TTCTACAAGATTCC

Indel alignment

Reference Sequence

TACAGCTGAACTGAACTCCTGCCTGTACAGCTCGTTTCTACAAGATTCCAGACCTGGAA

TGCCTGTACAGCTCGTT - TCTACAAGATT C

AACTGAACTCCTGCCTGTACAGCTCGT

TGAAACTCCTGCCTGTACAGCTCGTT - TCTA

TGTACAGCTCGTT - TCTACAAGATTCCAGA

CTCCTGCCTGTACAGCTCGTT - TCTACAAG

ACTCCTGCCTGTACAGCTCGTT - TCTACAA

GCCTGTACAGCTCGTT - TCTACAAGATTCC

Variant Calling

Reference Sequence

TACAGCTGAACTGAACTCCTGCCTGTACAGCTCGTTTCTACAAGATTCCAGACCTGGAA

TGCCTGTACAGCTCGTTTCTACAAGATTTC

AACTGAACTCCTGCCTGTAC**G**GCTCGT

TGAACTCCTGCCTGTAC**A**GCTCGTTTCTA

TGTAC**A**GCTCGTTTCTACAAGATTCCAGA

CTCCTGCCTGTAC**G**GCTCGTTTCTACAAG

ACTCCTGCCTGTACAGCTCGTTTCTACAA

GCCTGTAC**G**GCTCGTTTCTACAAGATTCC



Bayes Theorem

$$P(\text{Ref}|\text{Data}) = \frac{P(\text{Data}|\text{Ref}) * P(\text{Ref})}{P(\text{Data})}$$

- Incorporates the probability of our data given the hypothesis
- Incorporates prior information – the probability this site is reference.
- Incorporates the probability of the data under all hypotheses
- Provides a probability of belief in the hypothesis

Variant quality score

- The probability that the site was incorrectly called
- Phred scaled so its more intuitive
- Can be used to filter low quality sites

Genotype quality score

- The probability that the site was incorrectly genotyped.
- Is similar to VQS when calling single individual's variants – but not when calling population variants.
- Can be used to filter low quality sites on an individual basis when calling variants on a population.
- This is the value that you want to use to evaluate the quality of an variant.

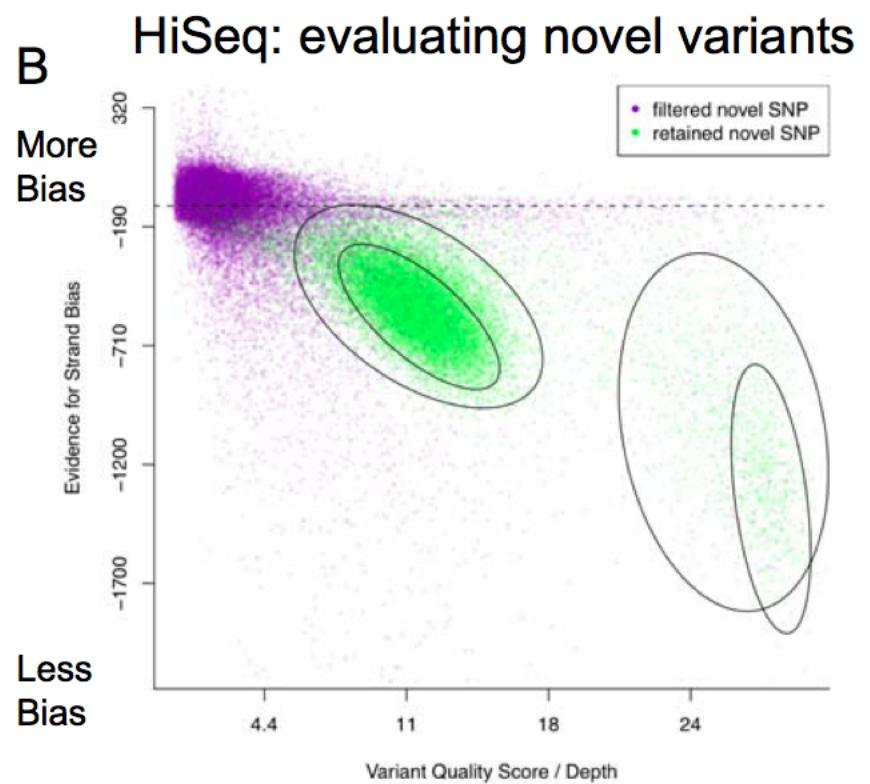
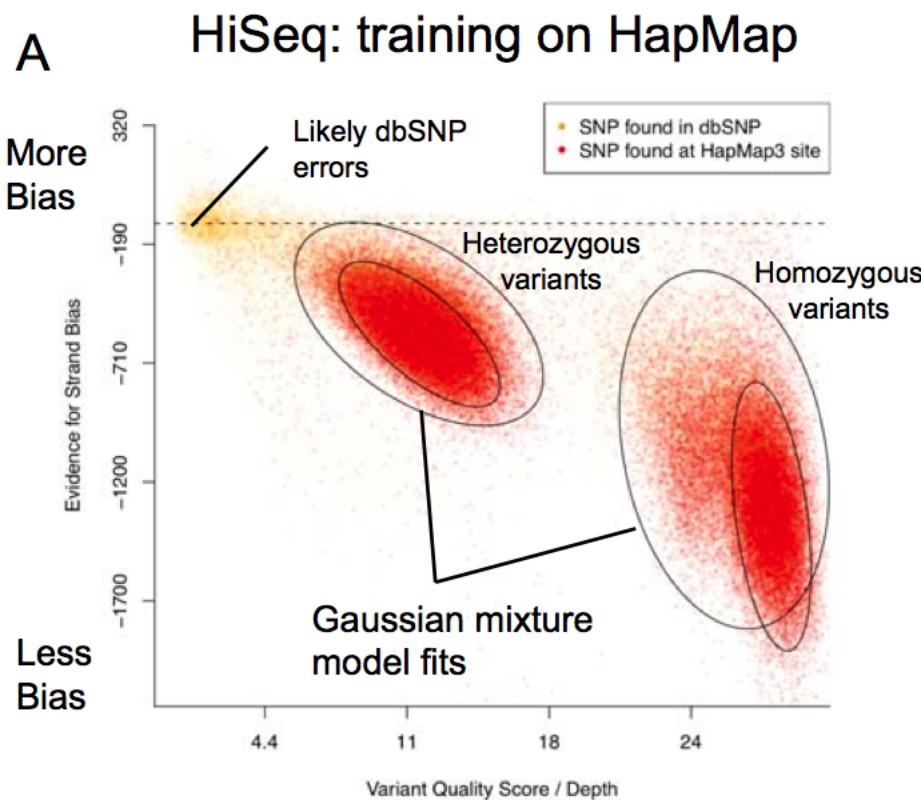
Variant filtering and recalibration

- Filtering on VQS or GQS is so yesterday!
- Population based variant calling allows population based covariant analysis of VQS
- Variant calling programs are implementing these recalibrations
 - GATK – VQSR
 - Real Time Genomics – AVR Score
 - FreeBayes

VQSR

- Considering only high-quality sites (HapMap3 concordant) build a mixture model considering various co-varying parameters
 - Allele Balance (ref/alt in hets)
 - Homopolymer run
 - MAPQ score
 - Strand Bias
 - Depth of coverage
- Apply that model to all variants until a given number of known sites have been recovered.

VQSR



Evaluating variant call sets

- Variant counts (~3.5-4M WGS, ~50K capture, ~20K coding, ~10K missense)
- Ti/Tv Ratio (~2.1-2.8)
- Concordance with known sites (HapMap3)
- NIST/GCAT – Tool for comparing pipelines
 - <http://www.bioplanet.com/gcat>