# The use of genomics to understand human disease

Jonathan Pevsner, Ph.D. Kennedy Krieger Institute

> October25, 2016 eScience IEEE 2016

# Outline

Introduction to genomics and human disease

Identifying a mutation causing a disease: Sturge-Weber

Genomic variation in autism spectrum disorder

# Definitions of bioinformatics and genomics

- Bioinformatics is the interface of molecular biology and computer science.
- It is the analysis of proteins, genes and genomes using computer algorithms and computer databases.
- Genomics is the analysis of genomes. The tools of bioinformatics are used to make sense of the quintillions of base pairs of DNA that are sequenced by genomics projects.



A genome is the collection of DNA that comprises an individual.The human genome is organized into 23 pairs of chromosomes (1-22, XX for girls, XY for boys).

Gene: Classically, a unit of hereditary information localized to a particular chromosome position and encoding one protein. It is a DNA sequence that makes RNA and that often then makes protein.





#### Second letter

		U	С	А	G		
	U	UUU UUC UUA UUA UUG	UCU UCC UCA UCG	UAU UAC UAA Stop UAG Stop	UGU UGC UGA UGG Trp	U C A G	
	С	CUU CUC CUA CUG	CCU CCC CCA CCG	CAU CAC His CAA CAA GIn	CGU CGC CGA CGG	U C A G	Thire
	A	AUU AUC AUA AUG Met	ACU ACC ACA ACG	AAU AAC AAA AAA AAG	AGU AGC AGA AGA AGG Arg	U C A G	letter
	G	GUU GUC GUA GUG	GCU GCC GCA GCG	GAU GAC GAA GAA GAG Glu	GGU GGC GGA GGG	U C A G	

First letter









# Genes are expressed at different times and places



# Time of development

Body region, physiology, pharmacology, pathology

15 () 9 81 81 81 **ار** 11 19 Г Х \*\* 21 ۲ 

# **Growth of GenBank**



Year





























http://humanorigins.si.edu/



#### Next-generation sequence technology: Illumina



ocuments and Settings\pevsner\Local Settings\Temp\XPgrpwise\bowtiefiles.txt] Project Tools Macros Window Help W. Document |∦ 🖻 🛍 으 ⊆ | 🎟 | № 🏥 🔏 🥻 | 🚾 🗟 🔊 | № 🚏 🔚 🔸 🗉 🕨 | ? 🔠 🖨 🔍 LA-CS:7:1:743:1919:GTTATAGAGAAAAATTTGATTTAAAT: 40 40 40 40 40 40 40 26 40 26 24 34 24 24 40 4 [A-CS:7:1:208:1926:GTATCATCGGCCATGGTCACTCATAT: 40 23 40 38 31 40 31 40 24 33 28 33 (A-CS:7:1:176:1936:GGGGGAAGTAATAGATTTACGGGTCA: 40 40 40 40 40 40 40 40 40 40 40 40 18 (A-CS:7:1:157:1959:GTTTCCTTATCTGTAGAAGGGGGTAA: 40 40) 40 40 40 40 40 40 40 40 38 40 40 40 40 40 40 40 40 (A-CS:7:1:876:1939:GCATTAGCAAACTTAAAAAAATGTTT: 40 40 40 40 40 40 40 A-CS:7:1:681:1981:GATTGAATATCAGGTCTGGTACAAAA: 40 38 40 40 40 37 40 40 40 40 34 35 33 A-CS:7:1:248:744: GTTGAAACTGTAAGTATATGGATACA: 40 40 23 40 19 40 40 9 29 19 14 27 20 36 9 40 [A-CS:7:1:625:1953:GAAACTATCTGTTTCTAGAGGCTTGT: 40 36 30 40 40 40 35 40 40 40 40 21 40 A-CS:7:1:650:1988:GAATTTTTCACCACCTTCTTTTCAAA: 40 40 31 40 40 40 40 40 40 40 40 32 39 A-CS:7:1:206:1844:GTGACACAGCTTGCAAAAGACTTTAA: 40 17 37 40 17 40 21 25 36 10 22 29 10 26

From Illumina:

raw sequence data includes short reads and quality scores

#### IGV view of the human genome (zoomed to 3 billion base pairs)

🚟 IGV		
File Genomes View Track	es Regions Tools GenomeSpace Help	
Human hg19		+
	1 3 5 7 9 11 13 15 17 19 21 X 2 4 6 8 10 12 14 16 18 20 22 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	Y
NA12874 exome Coverage	Zoom in to see coverage.	
NA12874 exome	Zoom in to see alignments.	
	sequence data for one individual	
NA12878 exome Coverage	Zoom in to see coverage.	
NA12878 exome	Zoom in to see alignments.	
	sequence data for another individual	
		-
RefSeq Genes	alered a section to a section of the	
	TOPOS	
6 tracks	458M of	898M

#### IGV view of one gene (zoomed to 300,000 base pairs)

📅 IGV											_				
File Genomes <u>V</u> iew Trad	<u>k</u> s	Regions Tools G	enomeSpace Help			ideogra	m o	f chr	omo	some	9				
Human hg19		chr9	chr	r9:80,333,191-80,648,2	219	lideogra						Ъ			
		p24.1 p2	2.3 p21.3 p21.	.1 p13.2 p11.2	q12	q13 q21.12 q21.2	q21.33	q22.32 q31	.1 q31.3 q	33.1 q33.3	q34.13				
												_			
			80,400	kЬ		314 Kb 80,500 kb			80,6	00 kb					
										1		-			
NA12874 exome Coverage		Zoom in to see coverage.													
NA12874 exome		Zoom in to see alignments.													
												-			
NA12878 exome Coverage						Zoom in to see coverage.									
NA12878 exome						Zoom in to see alignments,									
						-									
							( d)					▼ ▲			
RefSeq Genes						GNAQ									
							ех	cons	of a g	gene					
6 tracks chr	-9:8	0,480,953								48	6M of 898M	1			

#### IGV view of two exons (zoomed to 10,000 base pairs)



#### IGV view of one exon (zoomed to 1,000 base pairs)



#### IGV view of one exon (zoomed to 40 base pairs)

IGV	<b>D</b>	Tala	<u> </u>	0	u.L.												-		,						<u>_     ×</u>
Human hg 19	chr9	TOOIS	Genome	space	chr9:	80,412,4	473-80,4	12,51	2		Go	Ê	4	•	Ø [		¢ 🖓			-				11111	+
	p24	4.1 p	22.3	p21.3	p21.1	p1 3.2	p1	1.2		q12	q	13 q	21.12	q21.2	q2	1.33	q22.3	32 q	31.1	q31	.3 q	33.1	q33.3	q34.1	3
		1		80,412,4 	480 bp		I			80,412	:,490 Бр 	· 40 b	— qı	I		81	),412,50( 	) Бр			1		80,4	112,510 Бр I	
NA12874 exome Coverage	[0 - 1:29]																								
			,		- · ·			·				· .		· · ·			_			· ·					
NA12874 exome								-																	
NA12878 exome Coverage	[0 - 73]																								
NA12878 exome Junctions																									
	T G T G T G	A T A T A T	C C C C C C	C T C T C T	G G G		6 T 6 T 6 T	G	G G G G	G	A C A C	T T T	C G C G C G		C C C		T T T	A A A A A A	G G G	C C C	A C A C A C	AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA	T C T C T C		G G
NA12878 exome					G		G T G T G T G T	G G G		GGGG		T					T		G	C C	A C A C	A			GGG
Sequence 🔿	T G	A T	<u> </u>		6			G	<u> </u>	c G		+ i T			ř C		T		<u>c</u>	ř C		<u> </u>		$\frac{1}{1}$	
RefSeq Genes		I	G		T		T		P		V		R	1		R		L			V		D	Q	
am	ino	ac	ids								re	fe	re	nc	e	ge	nc	m	<b>le</b>	S	eq	ue	ene	ce	- -
6 tracks chr9	:80,412,494																						5	53M of 90	MOC
#### IGV view of one exon (zoomed to 60 base pairs)



## Human genome sequencing

We currently obtain whole genome sequences at 30x to 50x depth of coverage. For a typical individual:

- 2.8 billion base pairs are sequenced x 30
  = 100 billion base pairs of DNA
- 3-4 million single nucleotide variants (SNVs)
- ~600,000 insertions/deletions (indels)
- Cost (research basis) is < \$1500 per genome</p>
- We try to sequence mother/father/child trios

# Human genome sequencing: one purpose is to compare humans to animals

We want to understand what makes the human genome unique.We compare our genome to those of primates and other organisms across the tree of life.

This was a major goal of the Human Genome Project.



# Human genome sequencing: another purpose is to compare humans to each other

A second goal is to understand variation across human genomes.We compare genomes from different geographic (ethnic) groups. Currently we are in the process of sequencing >1 million genomes.

This is a major goal of the HapMap Project and the 1000 Genomes Project.

For Kennedy Krieger patients our goals are:

- improve diagnosis
- improve treatment
- offer genetic counseling (e.g. risk in siblings)

Genetic variation is responsible for the adaptive changes that underlie evolution.

Some changes improve the fitness of a species. Other changes are maladaptive and represent disease.

Medical perspective: pathological condition.

Molecular perspective: mutation and variation.

#### Projected global deaths (2002 to 2030)



http://www.who.int/whosis/whostat2007.pdf

# Four broad causes of disease phenotypes



This chart is not to scale, and all the categories are interconnected.A genomic disorder could be caused by a deletion in which loss of a single gene has a key role (e.g. RAII in Smith-Magenis syndrome)

Life is a relationship between molecules, not a property of any one molecule. So is therefore disease, which endangers life. While there are molecular diseases, there are no diseased molecules. At the level of the molecules we find only variations in structure and physicochemical properties. Likewise, at that level we rarely detect any criterion by virtue of which to place a given molecule "higher" or "lower" on the evolutionary scale. Human hemoglobin, although different to some extent from that of the horse, appears in no way more highly organized. Molecular disease and evolution are realities belonging to superior levels of biological integration. There they are found to be closely linked, with no sharp borderline between them. The mechanism of molecular disease represents one element of the mechanism of evolution. Even subjectively the two phenomena of disease and evolution may at times lead to identical experiences. The appearance of the concept of good and evil, interpreted by man as his painful expulsion from Paradise, was probably a molecular disease that turned out to be evolution. Subjectively, to evolve must most often have amounted to suffering from a disease. And these diseases were of course molecular.

#### Emile Zuckerkandl and Linus Pauling (1962)

Life is a relationship between molecules, not a property of any one molecule. So is therefore disease, which endangers life. While there are molecular diseases, there are no diseased molecules. At the level of the molecules we find only variations in structure and physicochemical properties. Likewise, at that level we rarely detect any criterion by virtue of which to place a given molecule "higher" or "lower" on the evolutionary scale. Human hemoglobin, although different to some extent from that of the horse, appears in no way more highly organized. Molecular disease and evolution are realities belonging to superior levels of biological integration. There they are found to be closely linked, with no sharp borderline between them. The mechanism of molecular disease represents one element of the mechanism of evolution. Even subjectively the two phenomena of disease and evolution may at times lead to identical experiences. The appearance of the concept of good and evil, interpreted by man as his painful expulsion from Paradise, was probably a molecular disease that turned out to be evolution. Subjectively, to evolve must most often have amounted to suffering from a disease. And these diseases were of course molecular.

Emile Zuckerkandl and Linus Pauling (1962)

## Outline

Introduction to genomics and human disease

Identifying a mutation causing a disease: Sturge-Weber

Genomic variation in autism spectrum disorder

# Sturge-Weber syndrome

A port-wine birthmark affects about 1:333 people. It varies in size and location.



# Sturge-Weber syndrome

A port-wine birthmark affects about 1:333 people. It varies in size and location.





Sturge-Weber syndrome affects < 1:20,000 people. It affects ~8% of individuals with a facial PW birthmark. Features of SWS can be highly variable, and may include:

- Port-wine birthmark (facial cutaneous vascular malformation)
- Seizures
- Intellectual disability
- Abnormal capillary venous vessels in the leptomeninges of the brain and choroid
- Glaucoma
- Stroke

## Sturge-Weber syndrome presentation



### Sturge-Weber syndrome presentation



SWS appears to be sporadic (rather than familial)

In some studies, identical twins are discordant (consistent with a model of somatic mosaicism)

Rudolf Happle (1987) speculated that a series of neurocutaneous disorders are caused by somatic mosaicism.

"A genetic concept is advanced to explain the origin of several sporadic syndromes characterized by a mosaic distribution of skin defects. It is postulated that these disorders are due to the action of a lethal gene surviving by mosaicism."



Somatic mosaic mutation

Somatic: changes occur in development (rather than being inherited).

Germline: perhaps an individual with such a mutation would not survive.

Mosaic: only part of the body is affected.



Fertilized egg (from which body's cells arise)

Fertilized egg divides, forms embryo

DNA in one cell becomes altered



G becomes A (in AKT1 or in GNAQ)

As the cells in the embryo divide, both normal and mutant cells expand and affect development

The baby's cells have normal or mutant gene Some parts of the body grow differently than those with normal cells

http://www.genome.gov/dmd/index.cfm?node=Photos/Graphics

Strategy: sequence and compare two genomes from each patient (n=3 individuals)



DNA from portwine birthmark (presumed affected)





Strategy: sequence and compare two genomes from each patient (n=3 individuals)



DNA from portwine birthmark (presumed affected)







DNA from blood (presumed unaffected)



Strategy: sequence and compare two genomes from each patient (n=3 individuals)

#### Each genome:

- ~3 billion bases of DNA
- Sequenced to 30x average depth of coverage, so 100 billion bases per genome
- A pair of genomes is compared (using a somatic variant caller)
- 100 GB raw data per genome
- Allow < I TB storage/genome</li>



#### The NEW ENGLAND JOURNAL of MEDICINE

ORIGINAL ARTICLE

### Sturge–Weber Syndrome and Port-Wine Stains Caused by Somatic Mutation in GNAQ

Matthew D. Shirley, Ph.D., Hao Tang, Ph.D., Carol J. Gallione, B.A., Joseph D. Baugher, Ph.D., Laurence P. Frelin, M.S., Bernard Cohen, M.D., Paula E. North, M.D., Ph.D., Douglas A. Marchuk, Ph.D., Anne M. Comi, M.D., and Jonathan Peysner, Ph.D.

#### ABSTRACT

#### BACKGROUND

The Sturge–Weber syndrome is a sporadic congenital neurocutaneous disorder characterized by a port-wine stain affecting the skin in the distribution of the ophthalmic branch of the trigeminal nerve, abnormal capillary venous vessels in the leptomeninges of the brain and choroid, glaucoma, seizures, stroke, and intellectual disability. It has been hypothesized that somatic mosaic mutations disrupting vascular development cause both the Sturge–Weber syndrome and port-wine stains, and the severity and extent of presentation are determined by the developmental time point at which the mutations occurred. To date, no such mutation has been identified.

# Analysis of high confidence results with Strelka resulted in one candidate mutation

Step	somatic SNVs	somatic indels
Pre-filtered calls	24,848	1,646
Post-filtered calls	1,300	27
VAAST pre-filtered	83	NA
VAAST post-filtered	1	NA

All 27 somatic indels were in repetitive regions

We performed targeted sequencing of a portion of GNAQ.

In skin samples, almost all patients had the mutation.

The mutant allele frequency was 1% to about 18%.

Patient No.	Mutation Present†	PWS	sws	Mutant Allele Frequency:	Total No. of Samples Assayed
				percent	
1	Yes	Yes	Yes	3.60	1
1	No	No	Yes	0.11	1
2	Yes	Yes	Yes	3.17	1
2	No	No	Yes	0.13	1
3	Yes	Yes	Yes	6.06-6.46	2
3	No	No	Yes	0.62-0.93	2
4	Yes	Yes	Yes	3.50-4.51	2
4	No	No	Yes	0.13-0.90	2
5	Yes	Yes	Yes	3.38	1
5	No	No	Yes	0.11	1
6	Yes	Yes	Yes	3.99	1
7	Yes	Yes	Yes	2.05-2.16	2
7	Yes	No	Yes	0.09-2.00	2
8	Yes	Yes	Yes	4.08	1
8	No	No	Yes	0.06	1
9	Yes	Yes	No	5.58	1
10	Yes	Yes	No	2.76	1
10	Yes	No	No	1.14	1
11	Yes	Yes	No	6.70	1
12	No	Yes	No	0.00	1
13	Yes	Yes	No	5.90	1
14	Yes	Yes	No	6.20	1
15	Yes	Yes	No	14.20	1
16	Yes	Yes	No	1.70	1
17	Yes	Yes	No	4.50	1
18	Yes	Yes	No	5.30	1
19	Yes	Yes	No	4.70	1
20	Yes	Yes	No	4.30	1
21	Yes	Yes	No	18.10	1
22	Yes	Yes	Yes	5.00	1

Table 1 Comptic Mutation of CNAO in Skin Complex

We performed targeted sequencing of a portion of GNAQ.

In skin samples, almost all patients had the mutation.

The mutant allele frequency was 1% to about 18%.

Patient No.	Mutation Present†	PWS	sws	Mutant Allele Frequency:	Total No. of Samples Assayed	
				percent		
1	Yes	Yes	Yes	3.60	1	
1	No	No	Yes	0.11	1	
2	Yes	Yes	Yes	3.17	1	
2	No	No	Yes	0.13	1	
3	Yes	Yes	Yes	6.06-6.46	2	
3	No	No	Yes	0.62-0.93	2	
4	Yes	Yes	Yes	3.50-4.51	2	
4	No	No	Yes	0.13-0.90	2	
5	Yes	Yes	Yes	3.38	1	
5	No	No	Yes	0.11	1	
6	Yes	Yes	Yes	3.99	1	
7	Yes	Yes	Yes	2.05-2.16	2	
7	Yes	No	Yes	0.09-2.00	2	
8	Yes	Yes	Yes	4.08	1	
8	No	No	Yes	0.06	1	
9	Yes	Yes	No	5.58	1	
10	Yes	Yes	No	2.76	1	
10	Yes	No	No	1.14	1	
11	Yes	Yes	No	6.70	1	
12	No	Yes	No	0.00	1	
13	Yes	Yes	No	5.90	1	
14	Yes	Yes	No	6.20	1	
15	Yes	Yes	No	14.20	1	
16	Yes	Yes	No	1.70	1	
17	Yes	Yes	No	4.50	1	
18	Yes	Yes	No	5.30	1	
19	Yes	Yes	No	4.70	1	
20	Yes	Yes	No	4.30	1	
21	Yes	Yes	No	18.10	1	
22	Yes	Yes	Yes	5.00	1	

### In brain samples, most (not all) patients had a mutation.

# Control brain samples: no mutation

Table 2. Somatic Mutation of GNAQ in Brain-Tissue Samples.\*

	Patient No.	Mutation Present	sws	Mutant Allele Frequency	Total No. of Samples Assayed
			_	percent	
	7	Yes	Yes	5.57-5.63	2
	23	Yes	Yes	5.56-5.78	2
	24	Yes	Yes	2.67-3.51	2
	25	No	Yes	0.02-0.10	2
	26	Yes	Yes	0.13-3.06	4
	27	Yes	Yes	2.19-5.12	2
	28	Yes	Yes	6.95-8.13	4
	29	Yes	Yes	6.04–11.15	5
	30	Yes	Yes	4.14	1
	31	Yes	Yes	4.78	1
	32	Yes	Yes	0.22-1.48	4
	33	Yes	Yes	4.04-5.74	2
	34	No	Yes	0.05-0.12	2
	35	Yes	Yes	0.05-1.51	7
	36	Yes	Yes	0.35-6.03	5
	37	Yes	Yes	5.74-6.49	2
	38	No	Yes	0.03-0.05	2
	39	Yes	Yes	1.83	1
	40	No	No	0.11	1
	41	No	No	0.05	1
	42	No	No	0.08	1
	43	No	No	0.09	1
	44	No	No	0.04	1
	45	No	No	0.04	1
	46	No	No, CCM	0.00	1
	47	No	No, CCM	0.00	1
	48	No	No, CCM	0.00	1
	49	No	No, CCM	0.00	1

# Targeted sequencing of a portion of GNAQ reveals mutations in SWS and PWS cases

# subjects	Tissue	SWS	GNAQ c.548 G->A	Detection
9	PWS	Yes	100%	Amplicon seq
7	Skin (non PWS)	Yes	14%	Amplicon seq
13	PWS	No	92%	Amplicon seq Primer extension
18	Brain	Yes	88%	Amplicon seq
6	Brain	No	0%	Amplicon seq
4	Brain	No: CCM	0%	Primer extension
669	Blood/LCL	N/A	0.7%	Exome seq

Amplicon sequencing: 13,000x median read depth Exome sequencing (IKG project): 271x median read depth Primer extension: SNaPshot assay (Doug Marchuk's lab)



- 13,000 reads
- Q30 base quality score
- I:1000 error rate
- Expect 13 errors in 13,000 reads
- If we see I 0x the error rate, call a mutation
- Call mutation if we see 130 T bases per 13,000 normal bases



# R183Q: an activating mutation in $G\alpha_q$

- In 2009 this identical mutation was described in uveal melanoma (a cancer involving melanocytes)
- The R183Q mutation occurs in 2-6% of these melanomas
- Another activating mutation (Q209L in Gαq) occurs in ~50% of uveal melanoma
- The mutation has been implicated in dermal hyperpigmentation



2007 Dorsam and Gutkind



2007 Dorsam and Gutkind



Mutations in genes encoding many of these signaling proteins<br/>cause somatic, mosaic, and often neurocutaneous disorders.TSC1, TSC2: tuberous sclerosisGNAQ: Sturge-WeberNF1: neurofibromatosisGNAS: McCune-AlbrightAKT1: Proteus syndromeRAS: epidermal neviPI3K: CLOVE syndrome, hemimegalencephalyNote the state of the state


We identified mutations in the GNAQ gene as the main cause of Sturge-Weber syndrome and port-wine birthmarks.

Knowing the genetic cause of the disease offers us a direction to search for treatments (and cures).

The consequence of the GNAQ mutation is to activate a cellular pathway. We can test drugs for the ability to reduce this persistent activation.

The same strategies may apply to treating uveal melanoma.

#### Outline

Introduction to genomics and human disease

Identifying a mutation causing a disease: Sturge-Weber

Genomic variation in autism spectrum disorder

- Deficits in social communication and interaction
- Restricted and repetitive patterns of behavior, interests or activities
- Symptoms cause significant impairment of function
- Diagnosed in childhood
- Comorbidities: intellectual disability, seizure, developmental delay, self-injury

Causes of ASD

- Associated with syndromic disorders (12% of ASD cases)
  - Fragile X syndrome
  - Rett Syndrome
  - Tuberous sclerosis
- de novo CNVs (6% of simplex cases)
- de novo SNVs/Indels (21% of simplex cases)

Heritability is the proportion of phenotypic variance due to genetic variance. For ASD, 50% to 90% heritability.













RESEARCH ARTICLE

#### The Contribution of Mosaic Variants to Autism Spectrum Disorder

Donald Freed<sup>1,2</sup>, Jonathan Pevsner<sup>1,2,3</sup>\*

#### Somatic mosaic variation in autism



#### Somatic mosaic variation in autism



# Collections of genotype and phenotype data from individuals with ASD

- Patients at the Kennedy Krieger Institute (50 trios)
- Simons Simplex Collection (SSC)
- MSSNG Project

Large collections of genomic data (e.g. 10,000 genomes) are available to qualified researchers: "the democratization of science."

# Collections of genotype and phenotype data from individuals with ASD

- Patients at the Kennedy Krieger Institute (50 trios)
  Simons Simplex Collection (SSC)
- MSSNG Project

# The Simons Simplex Collection (SSC)

- 8,938 individuals
  - 2,388 probands
  - I,774 siblings
  - 4,776 parents
- Simplex autism diagnoses
- DNA purified from blood
- Whole-exome sequencing on an Illumina platform
- Aligned sequence data publicly available on NDAR / AWS

#### Methods overview: finding mosaic variants

#### GATK pipeline (Genome Analysis Toolkit)

- Variant calling
- Genotyping
- Variant Quality Score Recalibration
- Identification of de novo variants
- Variant effect annotation
- Identification of mosaic variants

# Variant calling approach: GATK haplotype caller



https://software.broadinstitute.org/gatk/documentation/article?id=4148

# Methods:Variant calling via cloud computing

- Amazon EC2 + S3
- Virtual machines
- StarCluster (EC2 toolkit)



- Common bioinformatics tools (e.g. samtools)
- Python applications, R



#### Methods: Variant calling via cloud computing

AWS S3





AWS EC2

#### Methods: Variant calling via cloud computing







## Methods:Variant calling via cloud computing



#### Methods: Variant calling via cloud computing







### Methods: Variant calling via cloud computing





AWS EC2

http://www.livescience.com/topics/dna-genes

#### Methods overview: finding mosaic variants

- GATK pipeline
  - Variant calling (ploidy 5)
  - Genotyping
  - Variant Quality Score Recalibration
- Identification of de novo variants
- Variant effect annotation
- Identification of mosaic variants

- Variants are called per sample (we want variant information across all samples)
- Joint genotyping assesses all samples in the cohort simultaneously
- Samples are re-assessed for the presence of variants



http://www.livescience.com/topics/dna-genes

AWS S3 PEVS AWS EC2

AWS S3 PEVS



AWS EC2

AWS S3 PEVS AWS EC2





AWS EC2

#### Methods overview: finding mosaic variants

- GATK pipeline
  - Variant calling (ploidy 5)
  - Genotyping
  - Variant Quality Score Recalibration
- Identification of de novo variants
- Variant effect annotation
- Identification of mosaic variants

#### Variant Quality Score Recalibration

- Variant calling and genotyping are subject to systematic biases
- False positive variants due to these biases can be identified and filtered
  - Machine learning (Gaussian mixture model)
  - Known true positive (and false positive) variants
  - Set sensitivity thresholds

#### Methods overview: finding mosaic variants

- GATK pipeline
  - Variant calling (ploidy 5)
  - Genotyping
  - Variant Quality Score Recalibration

Identification of de novo variants

- Variant effect annotation
- Identification of mosaic variants

#### Identification of De Novo Variants

- De novo variants are present in a child but not either parent
- Identified de novo variants using a hard-filter approach
  - Variant present in unrelated sample
  - Read depth (20x)
  - Minimum genotype quality (20)

#### Methods overview: finding mosaic variants

- GATK pipeline
  - Variant calling (ploidy 5)
  - Genotyping
  - Variant Quality Score Recalibration
- Identification of de novo variants

Variant effect annotation

Identification of mosaic variants

Type (alphabetical order)	Count	Percent
3_prime_UTR_variant	102	1.119%
5_prime_UTR_premature_start_codon_gain_variant	8	0.088%
5_prime_UTR_variant	57	0.625%
disruptive_inframe_deletion	34	0.373%
disruptive_inframe_insertion	6	0.066%
downstream_gene_variant	771	8.458%
frameshift_variant	213	2.337%
inframe_deletion	24	0.263%
inframe_insertion	6	0.066%
intergenic_region	115	1.262%
intragenic_variant	47	0.516%
intron_variant	2,356	25.845%
missense_variant	2,687	29.476%
non_coding_exon_variant	69	0.757%
protein_protein_contact	235	2.578%
splice_acceptor_variant	54	0.592%
splice_donor_variant	62	0.68%
splice_region_variant	364	3.993%
start_lost	4	0.044%
stop_gained	159	1.744%
stop_lost	5	0.055%
synonymous_variant	1,172	12.857%
upstream_gene_variant	566	6.209%

# Variant Effect Annotation

### Methods overview: finding mosaic variants

- GATK pipeline
  - Variant calling (ploidy 5)
  - Genotyping
  - Variant Quality Score Recalibration
- Identification of de novo variants
- Variant effect annotation

Identification of mosaic variants
## Identifying mosaic variants



https://www.genome.gov/imagegallery/

## Identifying mosaic variants



https://www.genome.gov/imagegallery/

Validating mosaic variants by phasing

We developed a workflow to identify high quality candidates from sequence data. We also developed methods to validate somatic variants by phasing.



haplotype 1 haplotype 2

Physical position

Validating mosaic variants by phasing

We developed a workflow to identify high quality candidates from sequence data. We also developed methods to validate somatic variants by phasing.



Physical position

# Identifying mosaic variants

- Binomial test
  - False discovery protection with FDR of 0.05
- Additional filters
  - Mosaic variants must be in Krumm or lossifov
  - Mosaic variants must have AARF of < 0.34</li>
- Callset metrics
  - I00% precision for variant presence
  - 85% precision for mosaic status

#### De novo calls: comparision two recent studies



## Analysis of mutation rates

- Compare probands and siblings within the same family
- Increased mutation burden indicates a "contributory" role in disease
  - Rate = number of mutations per exome
  - contributory rate = proband rate sibling rate
  - % contributory = contributory rate / proband rate
- Only mutations at 40x sites in the trio
- Rates normalized to the entire capture target

## Rates of germline de novo mutation



## Rates of germline de novo mutation



## Rates of germline de novo mutation



#### Rates of mosaic mutation



- Classified mosaic mutations are a mix of mosaic and germline de novo events
- Same for classified germline de novo
- What is the contribution of incorrectly classified variants?
- Model parameters
  - Errors in classification of mosaic status
  - Validation rates
  - Number of germline and mosaic mutations





- The contribution of classified mosaic variants is primarily due to mosaic variation
- Some contribution of classified germline variants comes from mosaic variation

33% of mosaic variants contribute to 5.1% of ASD cases

6% of germline variants contribute to 5.6% of ASD cases











## Conclusions

- We identified many mosaic mutations (221 of ~4000 de novo mutations, i.e. 5.4%).
- Mosaic mutations were significantly enriched in probands relative to siblings and contribute to ~5% of simplex autism spectrum disorder diagnoses.
- We did not detect mosaic variants in paired brain/heart samples, at our level of detection.
- Mosaic variation may contribute to other neuropsychiatric disorders.

<u>Pevsner lab</u> Matt Shirley (now at Novartis) Larry Frelin Donald Freed (graduate student)

Alexis Norris, Jeremy Thorpe, Ike Adeshina, Kyra Feuer, McKinzie Garrison

<u>Collaborators: Sturge-Weber syndrome</u> Anne Comi (KKI) Doug Marchuk and Hao Tang, Carol Gallione (Duke) Bernard Cohen (JHU) Paula North (Wisconsin)