



PACBIO®



# CRISPR/Cas9 Enrichment and Long-read WGS for Structural Variant Discovery

PacBio CoLab Session

October 20, 2017

# PACBIO SMRT SEQUENCING



## Sequel System

### Long Reads

average 10 to 15 kb

### High Consensus Accuracy

random errors produce QV50 consensus

### Uniform, Unbiased Coverage

no GC% or sequence complexity bias

### Epigenetic Characterization

simultaneous detection of DNA methylation

# APPLICATIONS OF SMRT SEQUENCING



**Sequel System**

*De novo* genome assembly

Full isoform sequencing

Epigenetic characterization

Minor variant discovery

Structural variant discovery

Targeted sequencing

# APPLICATIONS OF SMRT SEQUENCING



**Sequel System**

*De novo* genome assembly

Full isoform sequencing

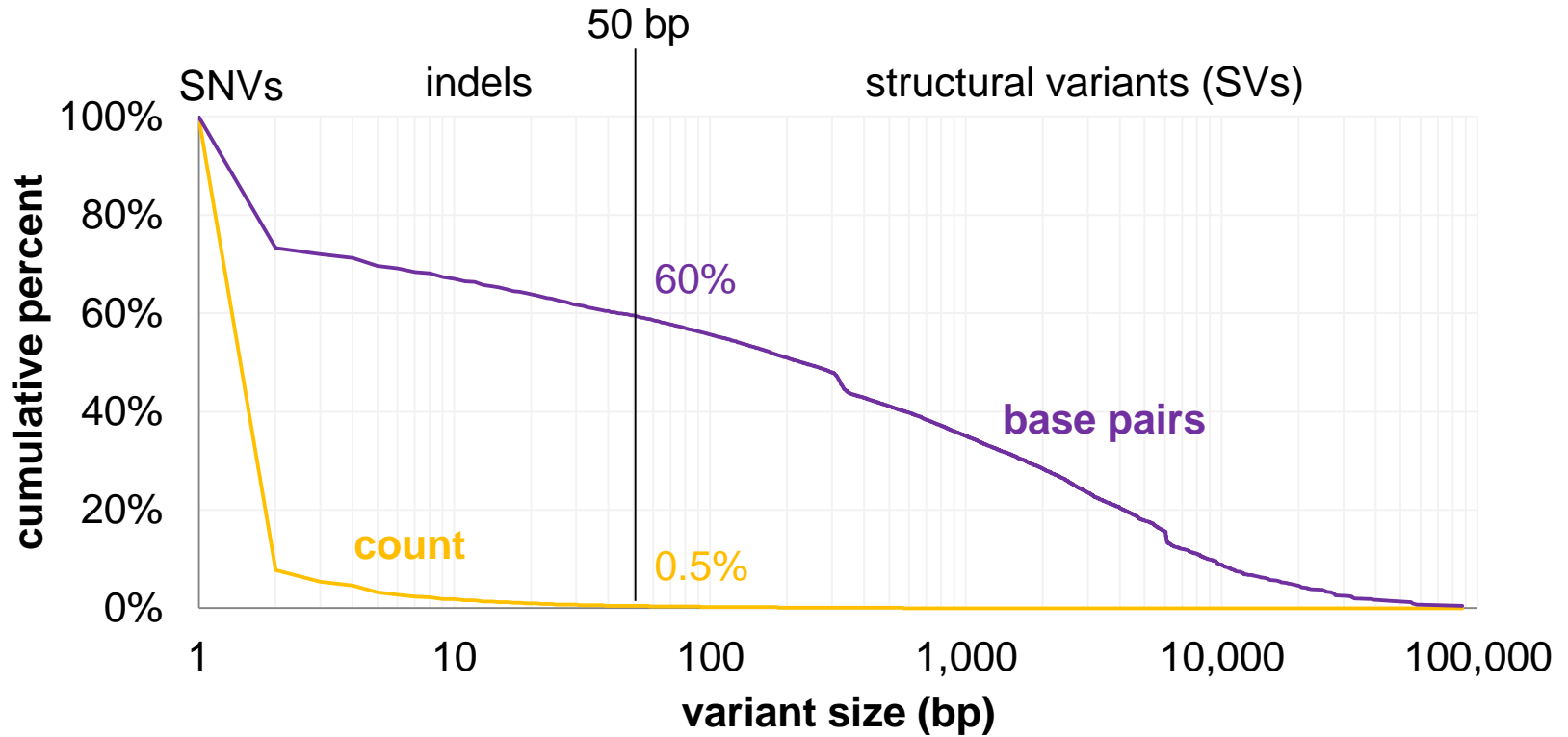
Epigenetic characterization

Minor variant discovery

Structural variant discovery

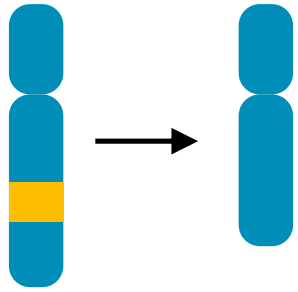
Targeted sequencing

# VARIATION IN A HUMAN GENOME – HG00733

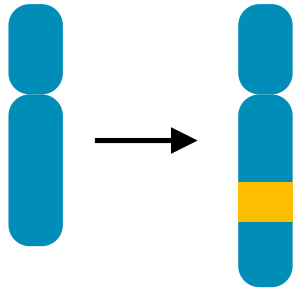


# TYPES OF STRUCTURAL VARIATION

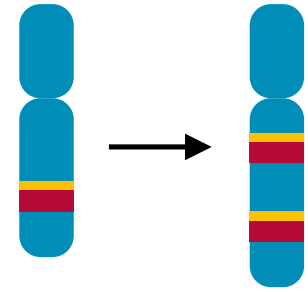
deletion



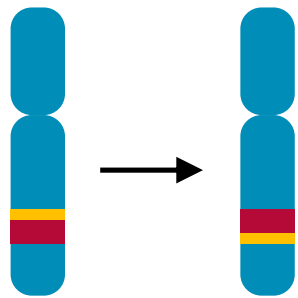
insertion



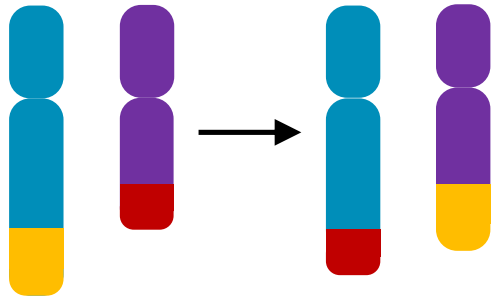
duplication



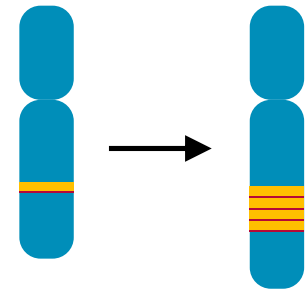
inversion



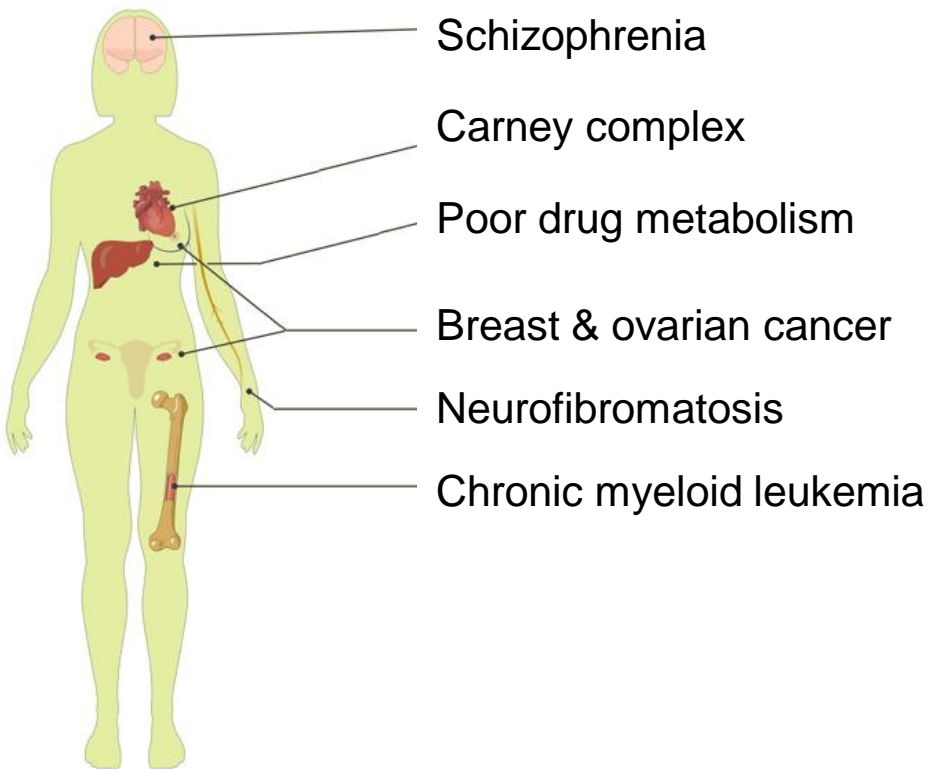
translocation



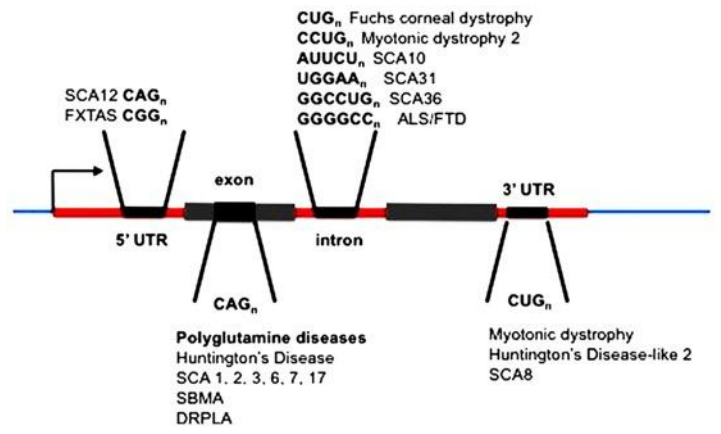
repeat expansion



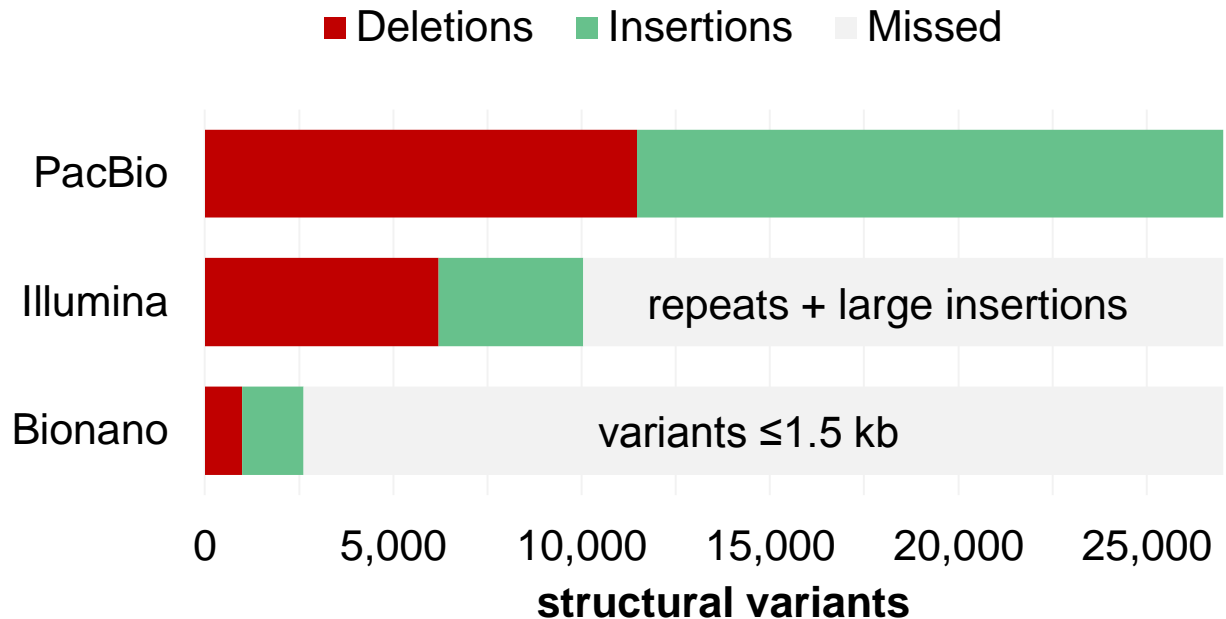
# STRUCTURAL VARIANTS AND DISEASE



## repeat expansion disorders



# TECHNOLOGY TO DETECT STRUCTURAL VARIANTS





“A move forward to **full-spectrum SV detection** ... will increase the diagnostic yield in patients with genetic disease, SV-mediated mutation, and repeat expansions.”

# **PacBio Long-Read WGS for Structural Variant Discovery**

**Targeted Enrichment without Amplification and SMRT Sequencing of Repeat-Expansion Disease Causative Genomic Regions**

# **PacBio Long-Read WGS for Structural Variant Discovery**

**Targeted Enrichment without Amplification and SMRT Sequencing of Repeat-Expansion Disease Causative Genomic Regions**

# FOR MORE INFORMATION – PACB.COM/SV



PRODUCTS + SERVICES

RESEARCH FOCUS

APPLICATIONS

SMRT SCIENCE

SUPPORT

COMPANY



## WHOLE GENOME SEQUENCING

- Human Whole Genome Sequencing
- Plant and Animal Whole Genome Sequencing
- Microbial Whole Genome Sequencing
- Structural Variation

## TARGETED SEQUENCING

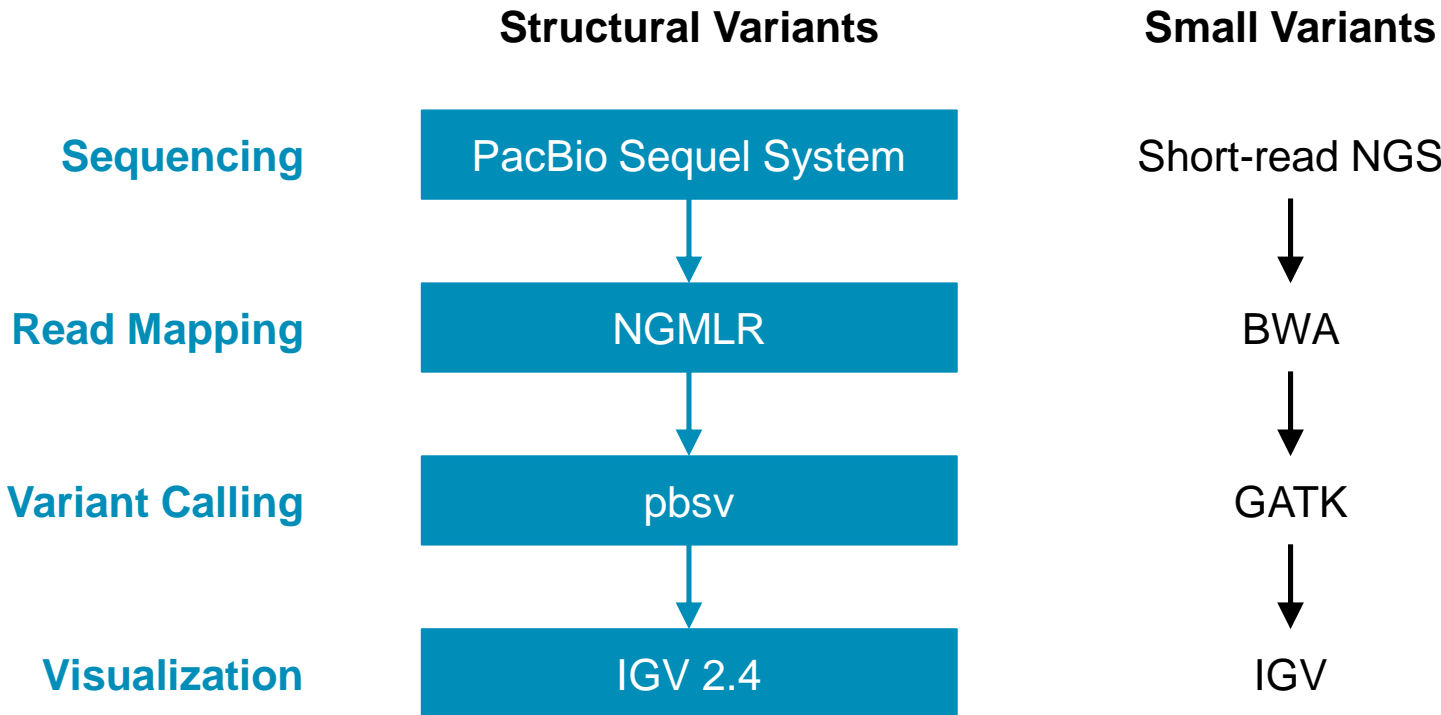
## CALLING ALL VARIANT TYPES

Structural variation accounts for most of the base pairs that differ between two human genomes, and causes many genetic disorders. The ability to study structural variants, in addition to smaller single nucleotide variants and indels, is critical to understanding how genetic variation impacts health and disease in the era of Precision Medicine.

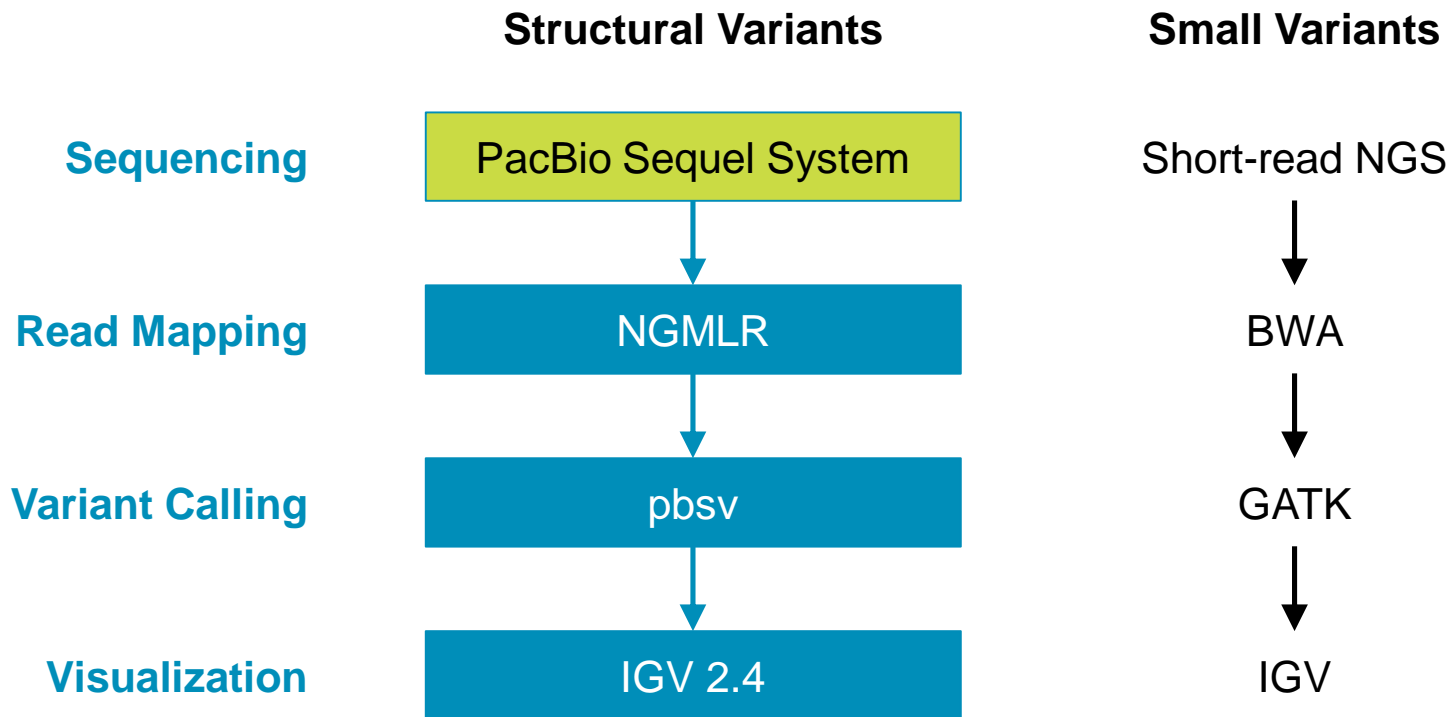
base pairs affected



# WGS FOR STRUCTURAL VARIANT DISCOVERY



# WGS FOR STRUCTURAL VARIANT DISCOVERY



# SEQUENCING

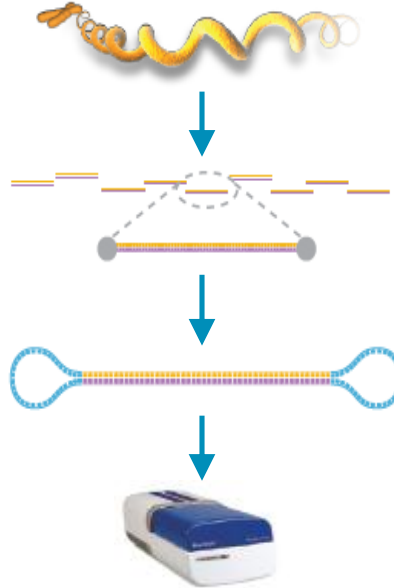
## Library Preparation

5  $\mu$ g DNA

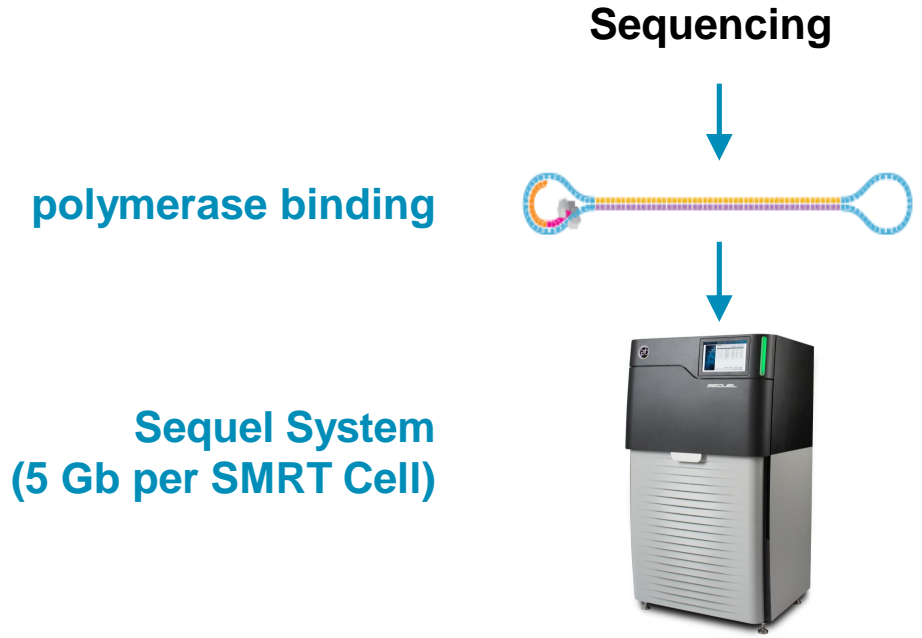
20 kb shear  
+ damage repair

SMRTbell adapter ligation

15 kb size selection

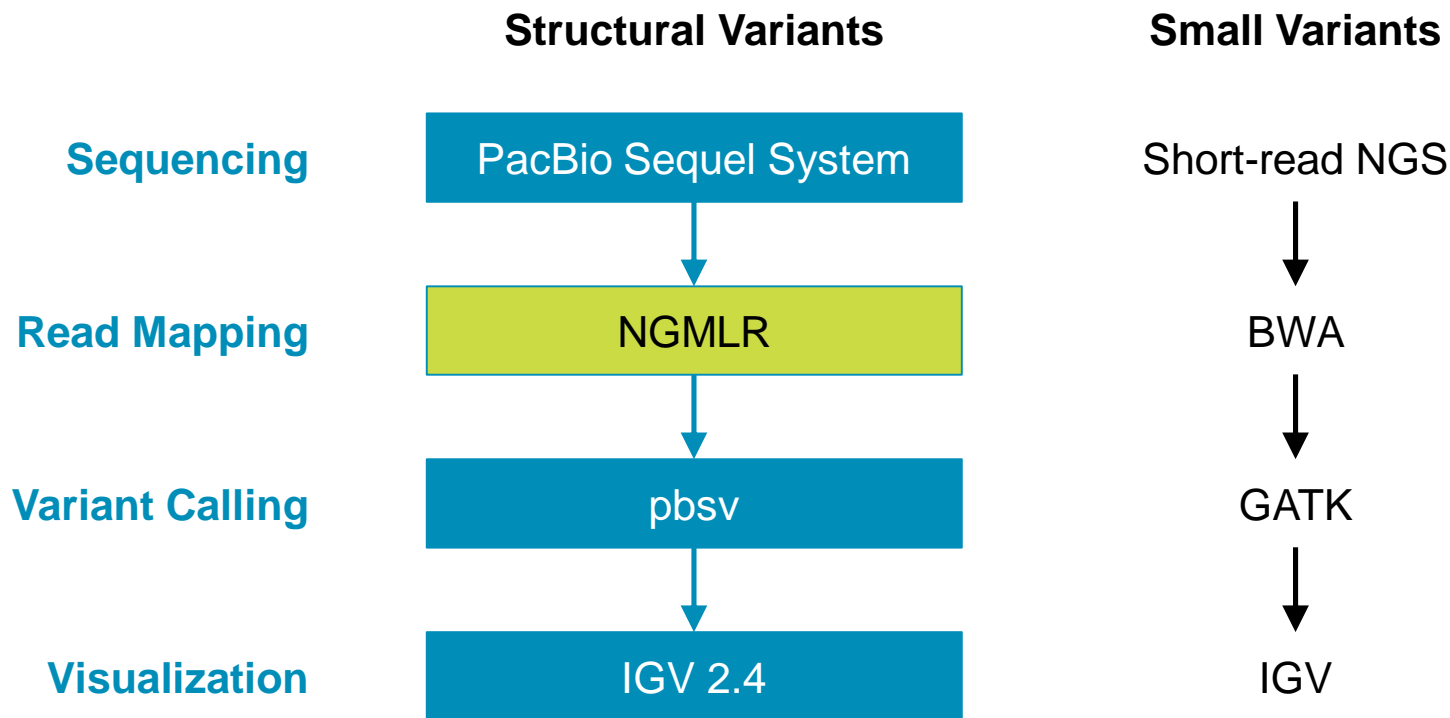


# SEQUENCING

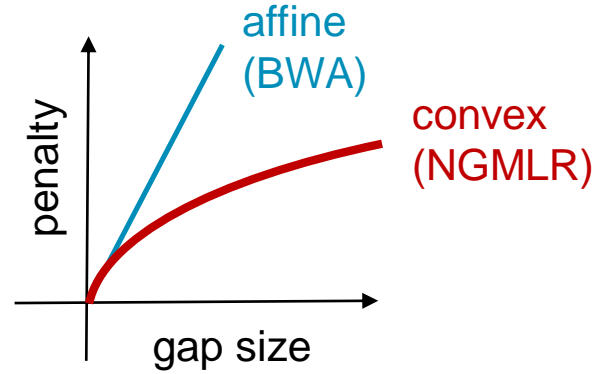
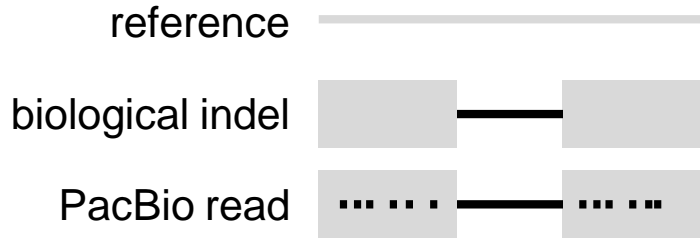




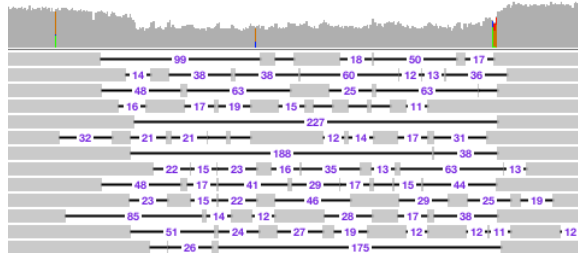
# WGS FOR STRUCTURAL VARIANT DISCOVERY



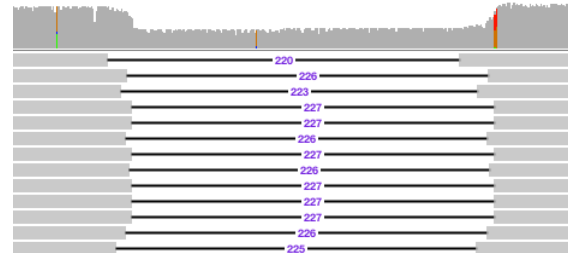
# READ MAPPING



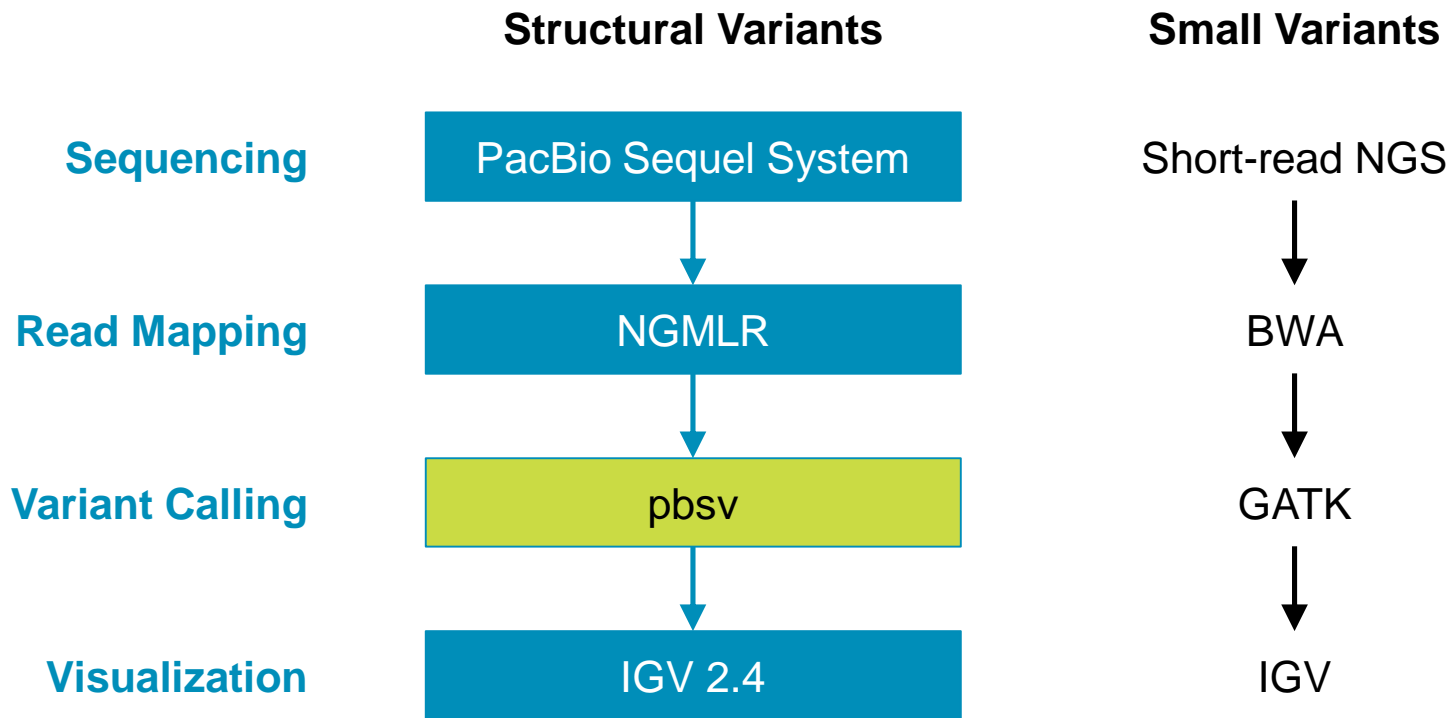
**BWA**



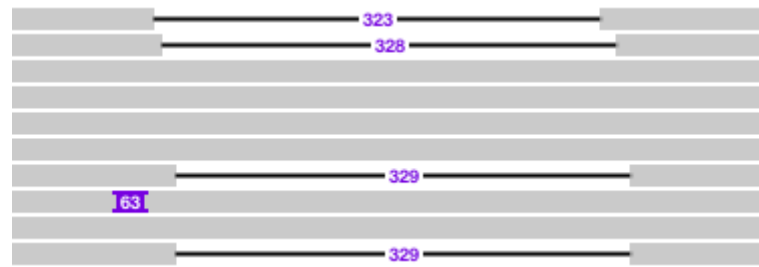
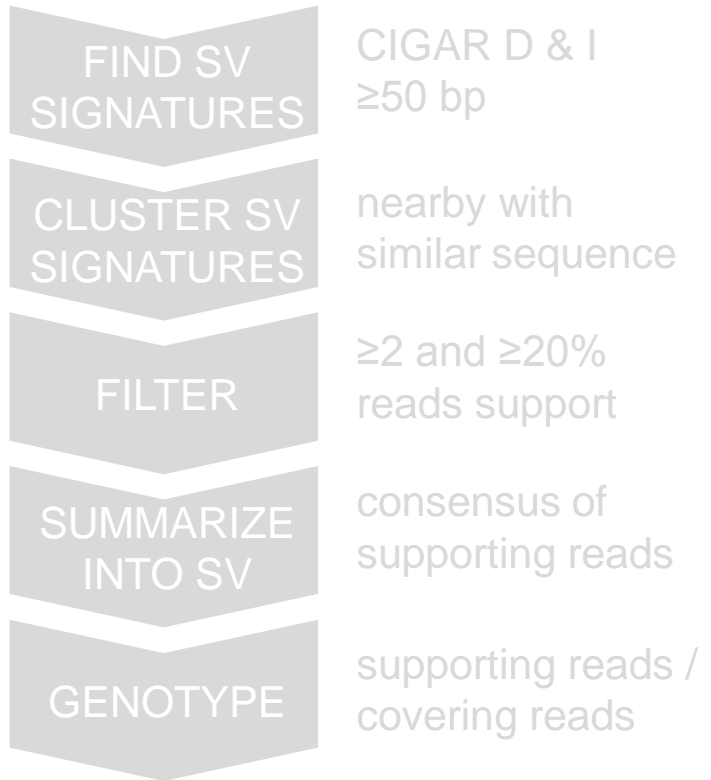
**NGMLR**



# WGS FOR STRUCTURAL VARIANT DISCOVERY



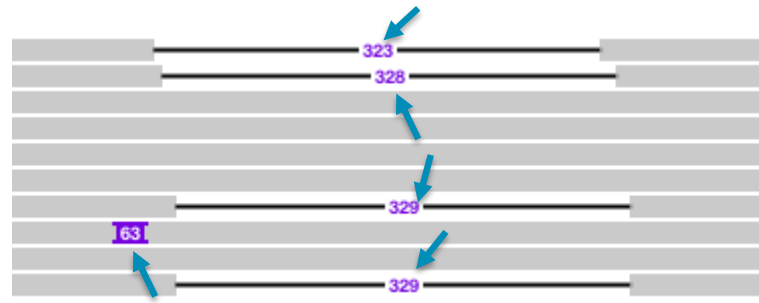
# VARIANT CALLING



# VARIANT CALLING



- CIGAR D & I  
≥50 bp
- nearly with similar sequence
- ≥2 and ≥20% reads support
- consensus of supporting reads
- supporting reads / covering reads



# VARIANT CALLING



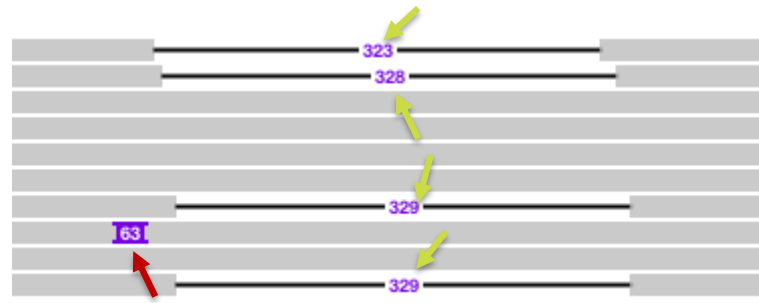
CIGAR D & I  
≥50 bp

nearby with  
similar sequence

≥2 and ≥20%  
reads support

consensus of  
supporting reads

supporting reads /  
covering reads



# VARIANT CALLING



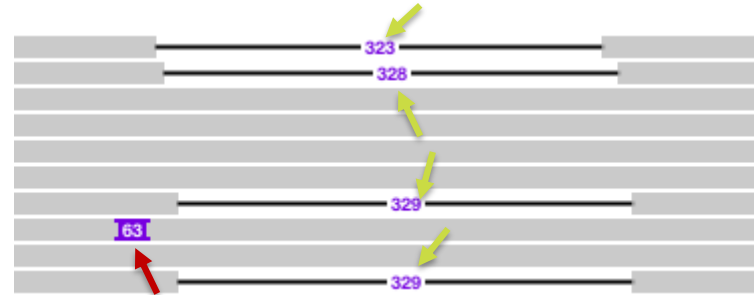
CIGAR D & I  
≥50 bp


nearly with  
similar sequence

≥2 and ≥20%  
reads support

consensus of  
supporting reads

supporting reads /  
covering reads



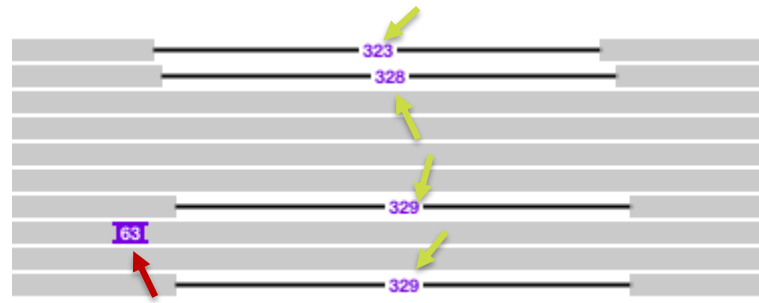
1 of 10  


4 of 10  


# VARIANT CALLING



- CIGAR D & I  $\geq 50$  bp
- nearby with similar sequence
- $\geq 2$  and  $\geq 20\%$  reads support
- consensus of supporting reads
- supporting reads / covering reads



1 of 10

4 of 10

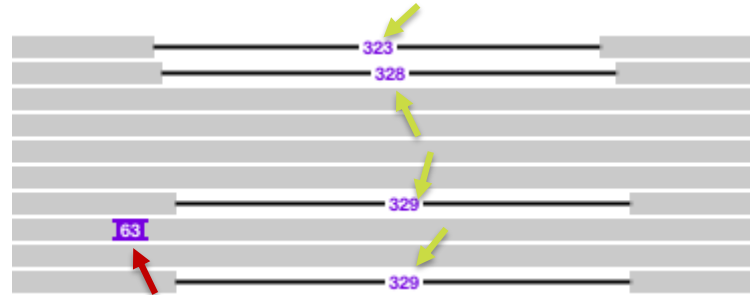
329 bp deletion



# VARIANT CALLING



- CIGAR D & I  $\geq 50$  bp
- nearly with similar sequence
- $\geq 2$  and  $\geq 20\%$  reads support
- consensus of supporting reads
- supporting reads / covering reads



1 of 10

4 of 10

329 bp deletion

heterozygous (4 of 10)

# VARIANT CALLING

**SMRT Analysis**

Create New Analysis - Settings [CANCEL] [START]

Name \*  
HG00733 10-fold SV

Analysis Application \*  
Structural Variant Calling [Beta]

Data Sets

Name
<input checked="" type="checkbox"/> HG00733_Subreads

References \*  
hg38

Structural Variants

minimum reads that support variant (count) ⓘ  
2

minimum reads that support variant (percent) ⓘ  
0.2

**SMRT Analysis**

Data

File Downloads	File	Size	Type
Analysis Log	Analysis Log	0 bytes	log
	Structural variants	12,426,862 bytes	vcf
	Structural variants	6,823,102 bytes	bed
	Aligned reads	18,542,771,404 bytes	bam
	Master Log	633 bytes	log

chr1  
904490  
ACGGGGCGGCTCC TCC TC CGA ACG TG GCC TCC TC CGA ACG CG GCC GGC TC CTC CTC CG AAC GCG GC CGC CTC CT OCT CCGA  
A  
PASS  
IMPRECISE;SVTY PE=DE L; END=90 458 7; SVLEN=-97; SVANN=T AN DEM  
GT:AD:DP  
0/1:9:15

**SMRT Analysis**

Report

Count by Annotation

Count by Annotation	Insertions (count)	Insertions (total bp)	Deletions (count)	Deletions (total bp)	All Variants (count)	All Variants (total bp)
Length Histogram (<1 Mb)						
Tandem Repeat	7,483	2,742,385	4,210	1,247,305	11,693	3,989,690
ASU	1,236	395,032	1,177	367,370	2,413	755,402
L1	44	244,741	83	444,635	127	689,376
SVK	18	31,987	29	51,831	47	83,818
Unannotated	4,344	2,007,459	2,861	2,803,452	7,205	4,810,911
Total	13,125	5,424,604	8,160	4,914,593	21,285	10,339,197

**SMRT Analysis**

Report

Length Histogram (<1 kb)

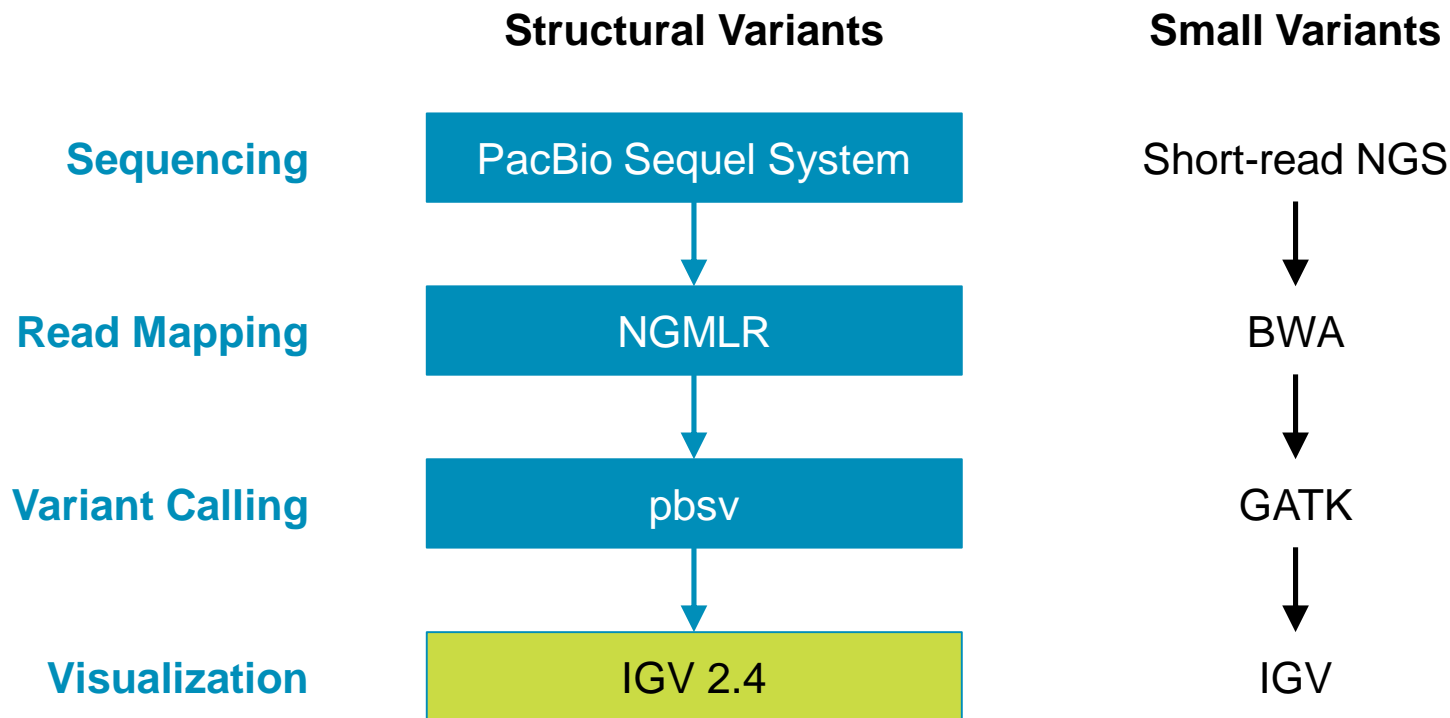
Count by Annotation

Length Histogram (<1 Mb)

Length Histogram (<1 Mb)

Variant Length (bp)	Deletions (count)	Insertions (count)
100	~3,000	~4,500
200	~1,500	~2,500
300	~1,000	~1,500
400	~500	~1,000
500	~200	~500
600	~100	~200
700	~50	~100
800	~20	~50
900	~10	~20
1,000	~5	~10

# WGS FOR STRUCTURAL VARIANT DISCOVERY



# VISUALIZATION

Home > Downloads

Downloads

Integrative Genomics Viewer - IGV 2.4

Install IGV

Note: IGV 2.4.x requires [Java 8](#).

Options for installing and running the current version of IGV:

1. Use the **Java Web Start** buttons below to launch IGV directly.
2. Download and run the **Mac application**; or
3. Download the **Windows zip archive** and run the `igv.bat` file using this version; or
4. Download the **binary distribution** and run IGV from the command line.

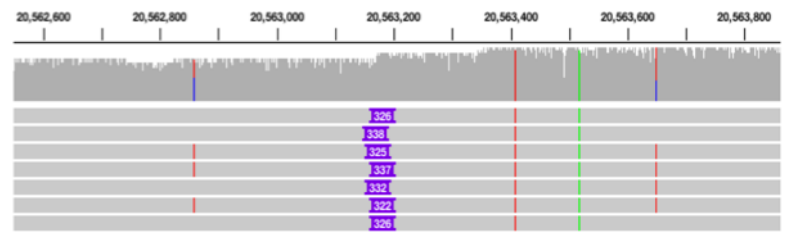
**1. Java Web Start.**

Clicking on one of the *Launch* buttons below will download a Java Web Start application.

- Some web browsers will only download the .jnlp file and not the application. You will need to manually download the application to your folder designated for browser downloads, and double-click it to launch.
- **JWS on Mac:** For some Mac users, IGV will not launch a notified of security errors, try the following instead:
  - Right-click on the downloaded .jnlp file and select `Open`.
  - Dismiss the warnings to continue.
  - When IGV has been run this way at least once from the command line, you can launch IGV normally.
  - Alternatively, install the **Mac App** or the **Binary Distribution**.
- **Windows:** To run with more than 1.2 GB of memory on Windows, you must **not include 64-bit Java by default, even if the operating system is 64-bit.** launch options with 32-bit Java will result in the error "could not find the main class".

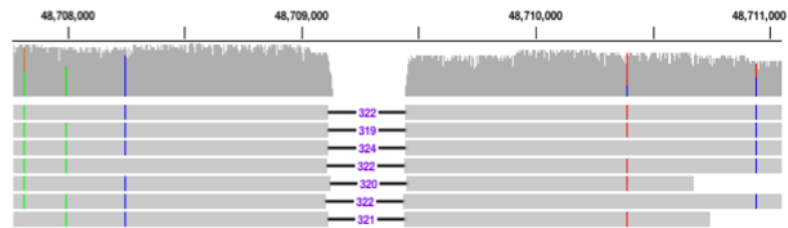
[Launch](#)      [Launch](#)

Launch with 750 MB      Launch with 1.2 GB (Max usable memory for Windows with 32-bit Java)



insertion

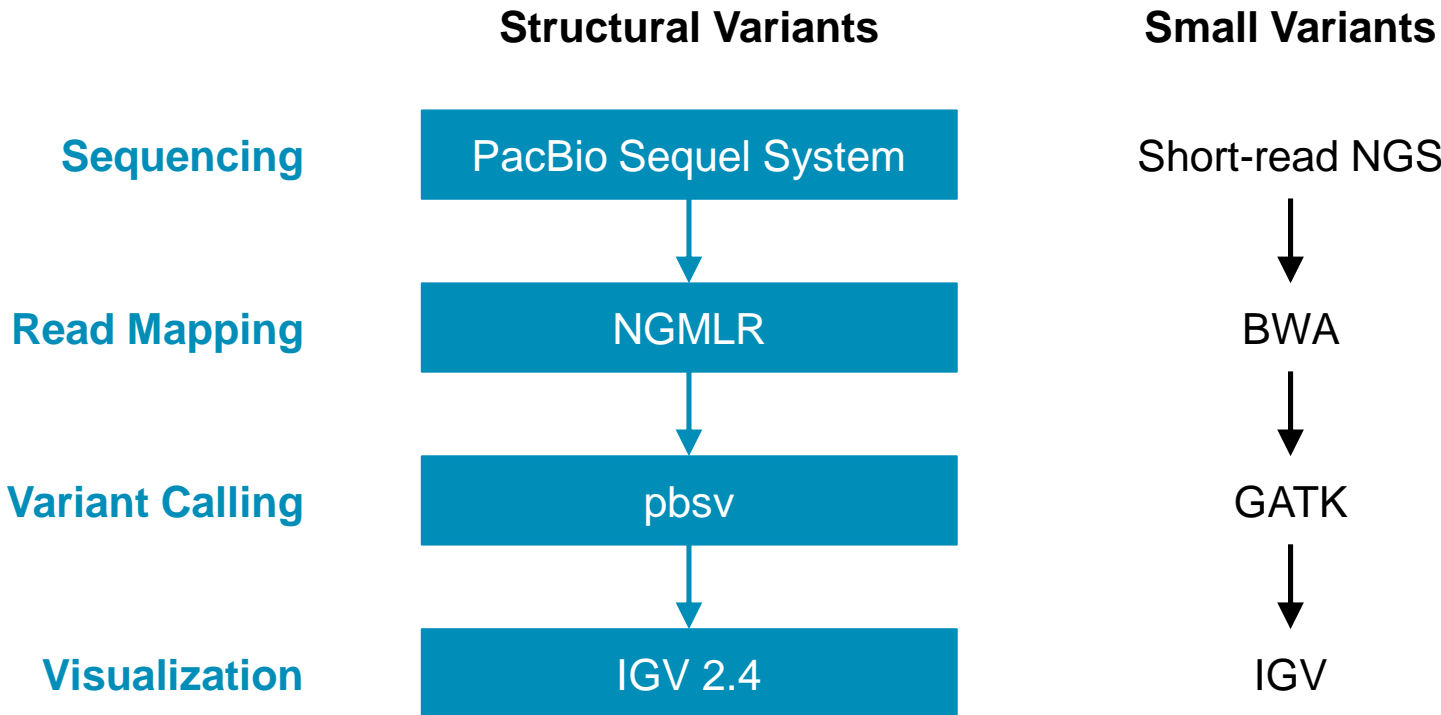
insertion label



deletion

deletion label

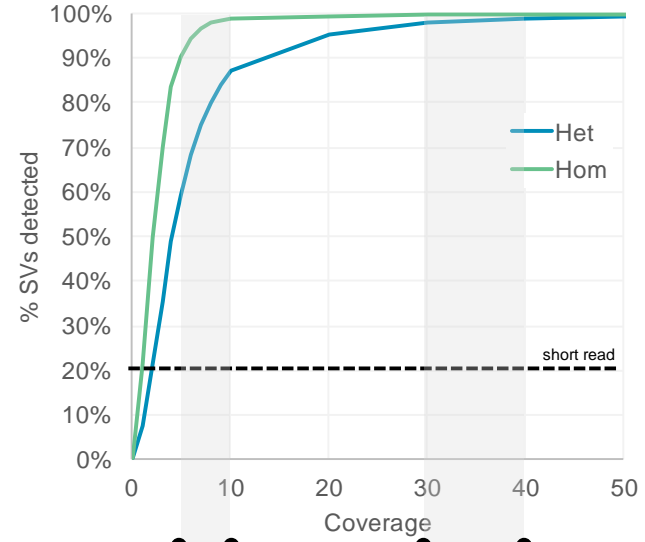
# WGS FOR STRUCTURAL VARIANT DISCOVERY



# HOW MUCH TO SEQUENCE?



Human HG00733  
Sequel System  
211 Gb (70-fold)



5- to 10-fold  
optimal tradeoff of  
cost vs. performance

disease gene discovery;  
population characterization

30- to 40-fold  
saturate discovery

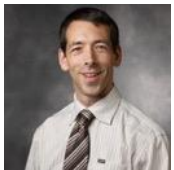
*de novo* variant discovery

# CLINICAL CASE HISTORY

- 7 yrs left atrial myxoma resection, atrial repair
- 10 yrs testicular mass, right orchiectomy
- 13 yrs pituitary tumor
- 16 yrs recurrence of myxomata, resection, adrenal microadenoma
- 18 yrs recurrence of ventricular myxomata, resection, VT
- 19 yrs ACTH-independent Cushing's disease, thyroid nodules
- 21 yrs transphenoidal resection of pituitary
- present (26 yrs) recurrence of myxomata, consideration for heart transplant

→ genetics suggests Carney complex  
PRKAR1A testing negative

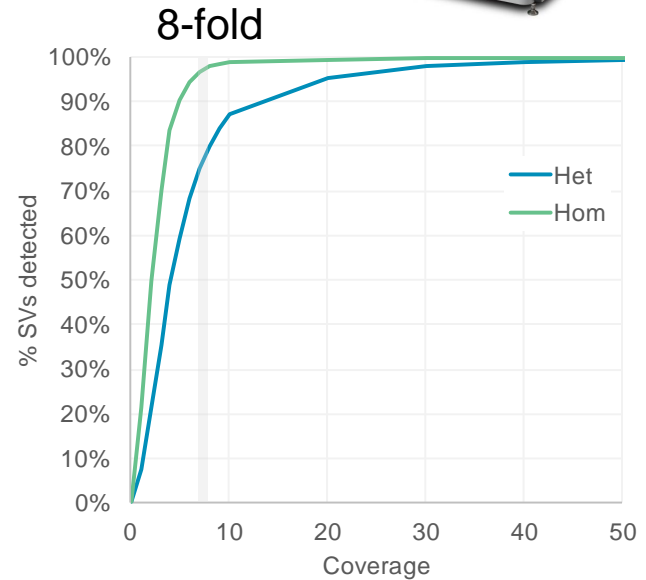
→ short-read whole genome sequencing negative



# EVALUATING STRUCTURAL VARIANTS

**Deletions    Insertions**

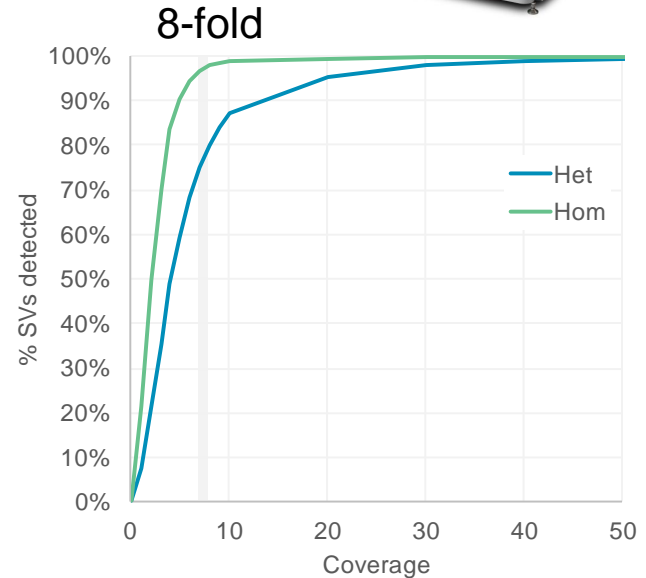
Initial call set    6,971    6,821





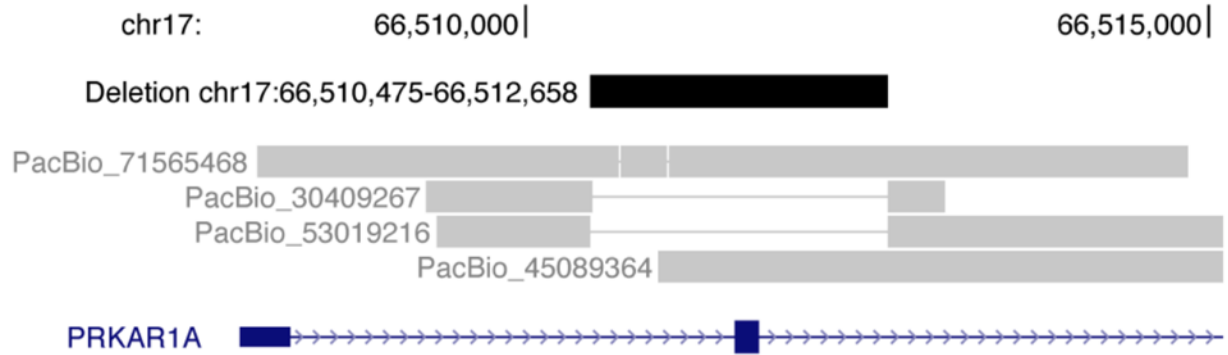
# EVALUATING STRUCTURAL VARIANTS

	Deletions	Insertions
Initial call set	6,971	6,821
Not in segdup	5,893	6,254
Not in NA12878 "healthy" control	2,476	3,171
Overlaps RefSeq coding exon	39	16
Gene linked to some disease in OMIM	3	3

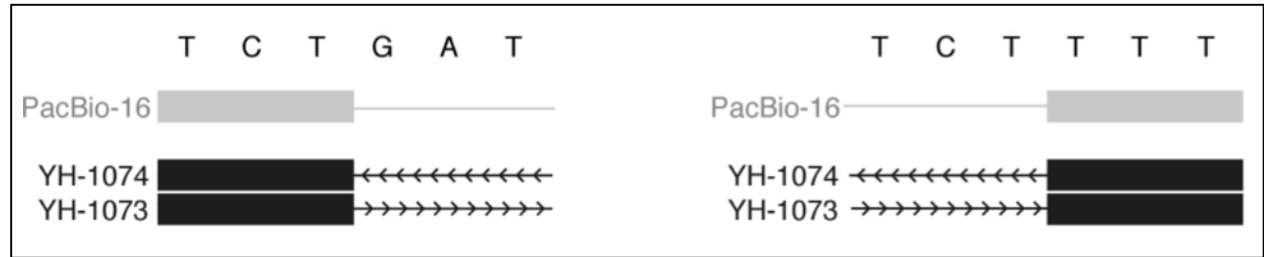


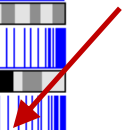
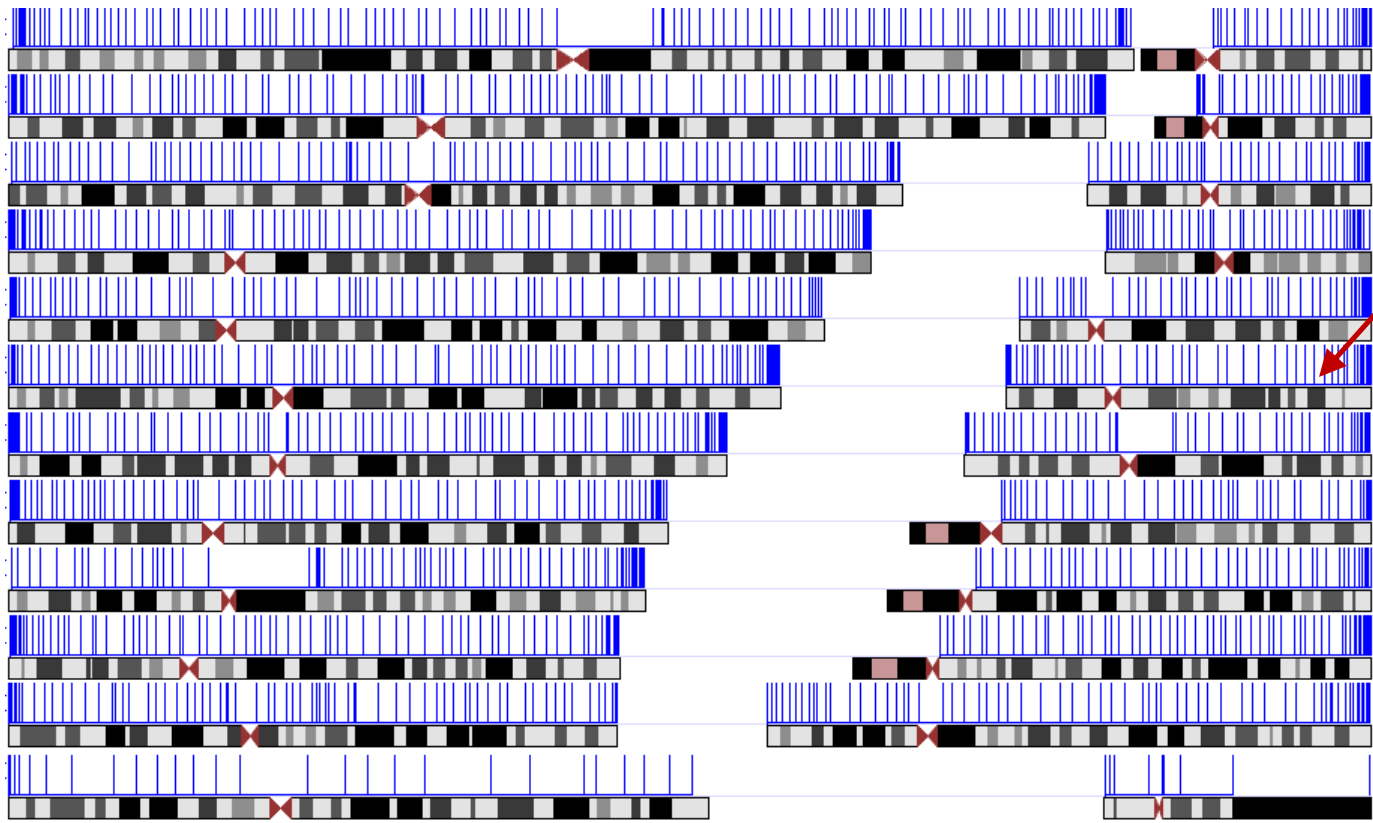
# HETEROZYGOUS 2.2 KB DELETION IN *PRKAR1A*

PacBio  
discovery



Sanger  
confirmation

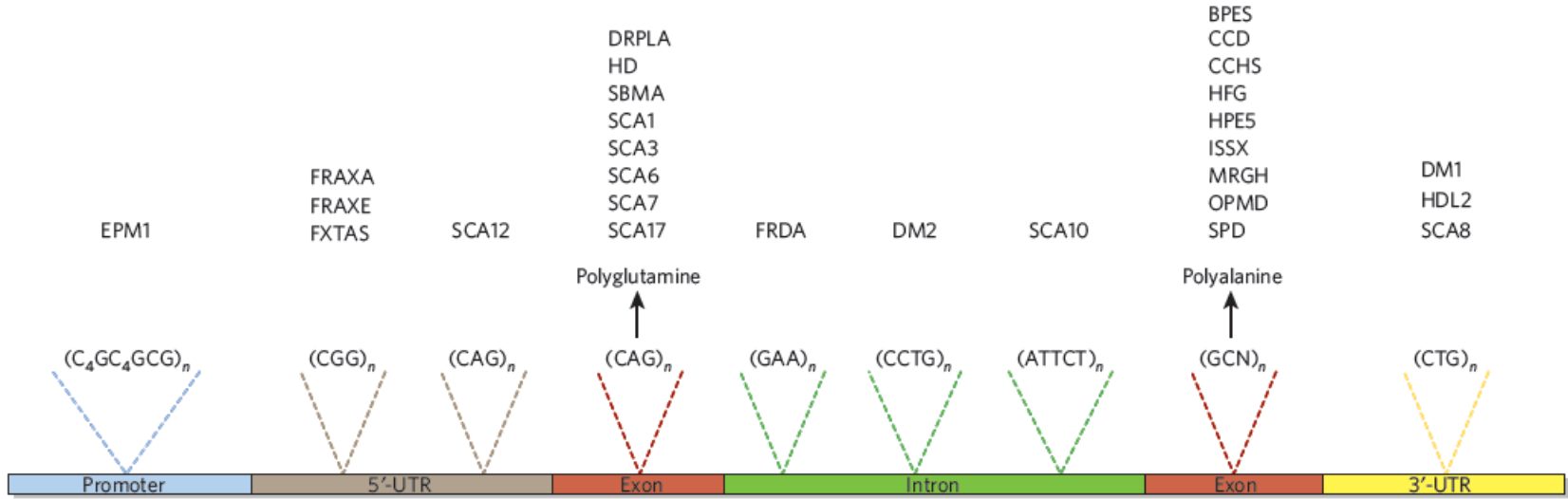




# **PacBio Long-Read WGS for Structural Variant Discovery**

**Targeted Enrichment without Amplification and SMRT Sequencing of Repeat-Expansion Disease Causative Genomic Regions**

# REPEAT EXPANSION DISEASES

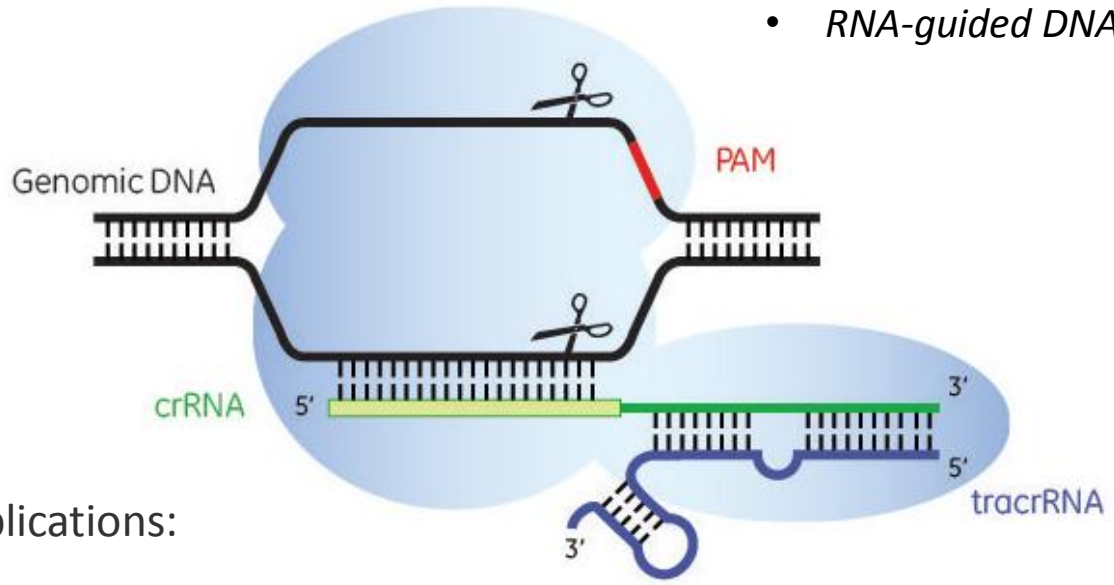


**Figure 1 | Location of expandable repeats responsible for human diseases.** The sequence and location within a generic gene of expandable repeats that cause human diseases are shown, and the associated diseases are listed. BPES, blepharophimosis, ptosis and epicanthus inversus; CCD, cleidocranial dysplasia; CCHS, congenital central hypoventilation syndrome; DM, myotonic dystrophy; DRPLA, dentatorubral-pallidolusian atrophy; EPM1, progressive myoclonic epilepsy 1; FRAXA, fragile X syndrome; FRAXE, fragile X mental retardation

associated with *FRAXE* site; FRDA, Friedrich's ataxia; FXTAS, fragile X tremor and ataxia syndrome; HD, Huntington's disease; HDL2, Huntington's-disease-like 2; HFG, hand-foot-genital syndrome; HPE5, holoprosencephaly 5; ISSX, X-linked infantile spasm syndrome; MRGH, mental retardation with isolated growth hormone deficiency; OPMD, oculopharyngeal muscular dystrophy; SBMA, spinal and bulbar muscular atrophy; SCA, spinocerebellar ataxia; SPD, synpolydactyly.

# CRISPR/CAS9 SYSTEM

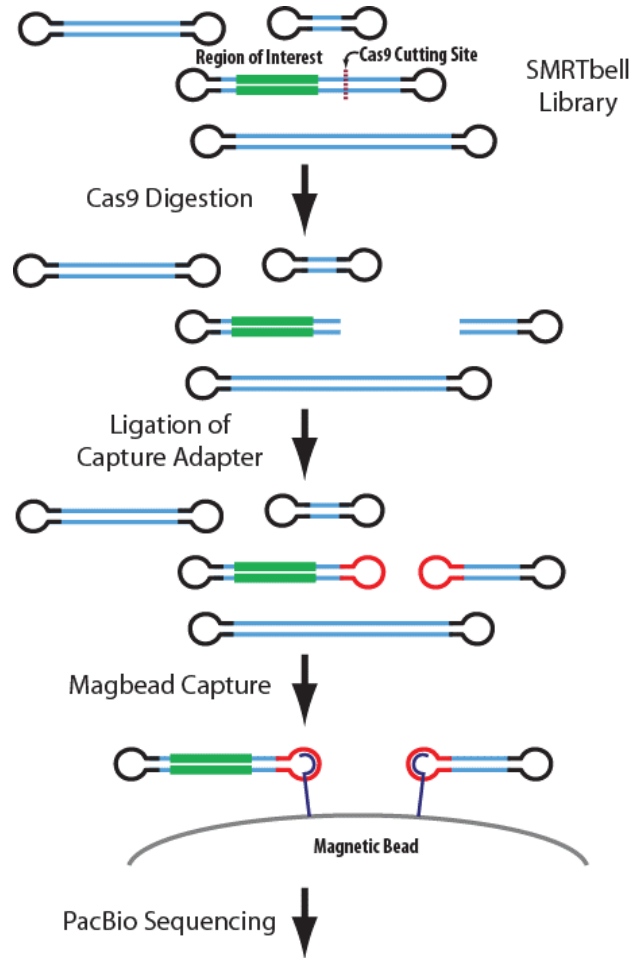
- *Bacterial Adaptive Immunity*
- *RNA-guided DNA Endonuclease*



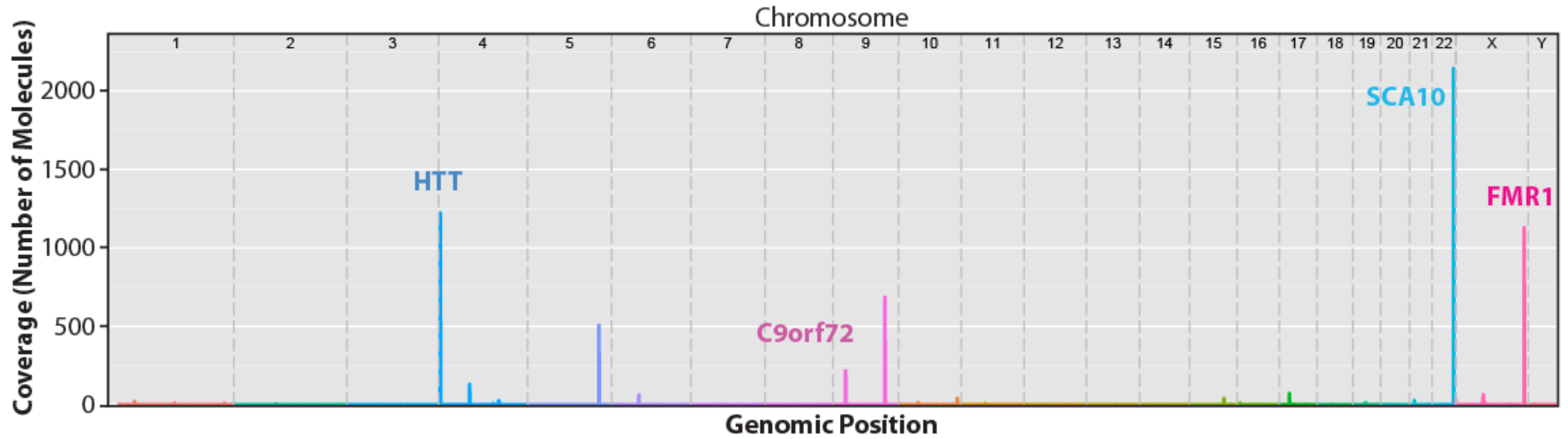
Some *in vivo* applications:

- Gene silencing
- Homology-directed repair
- Transient gene silencing or transcriptional repression
- Transient activation of endogenous genes
- Transgenic animals and embryonic stem cells

# PCR-FREE TARGET ENRICHMENT VIA CAS9



# COVERAGE ACROSS THE GENOME

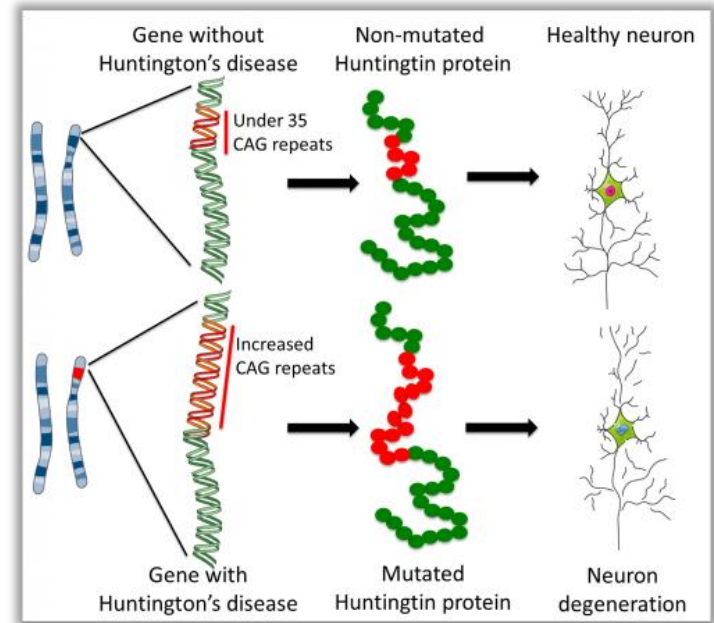
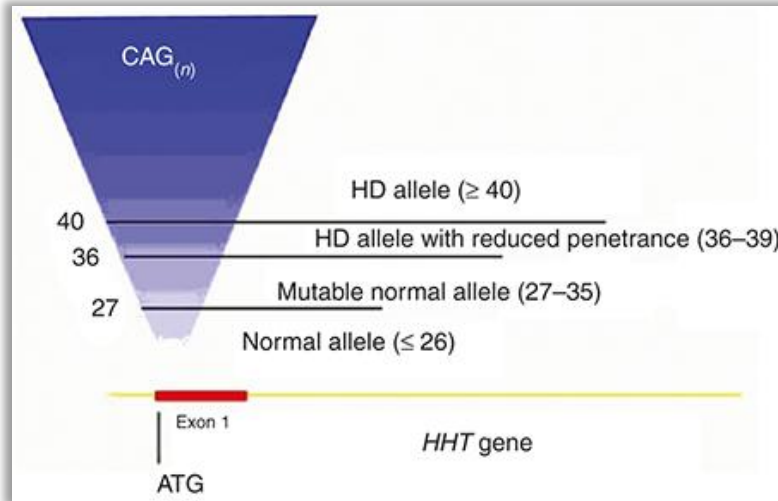


1 SMRT Cell (PacBio RS II)



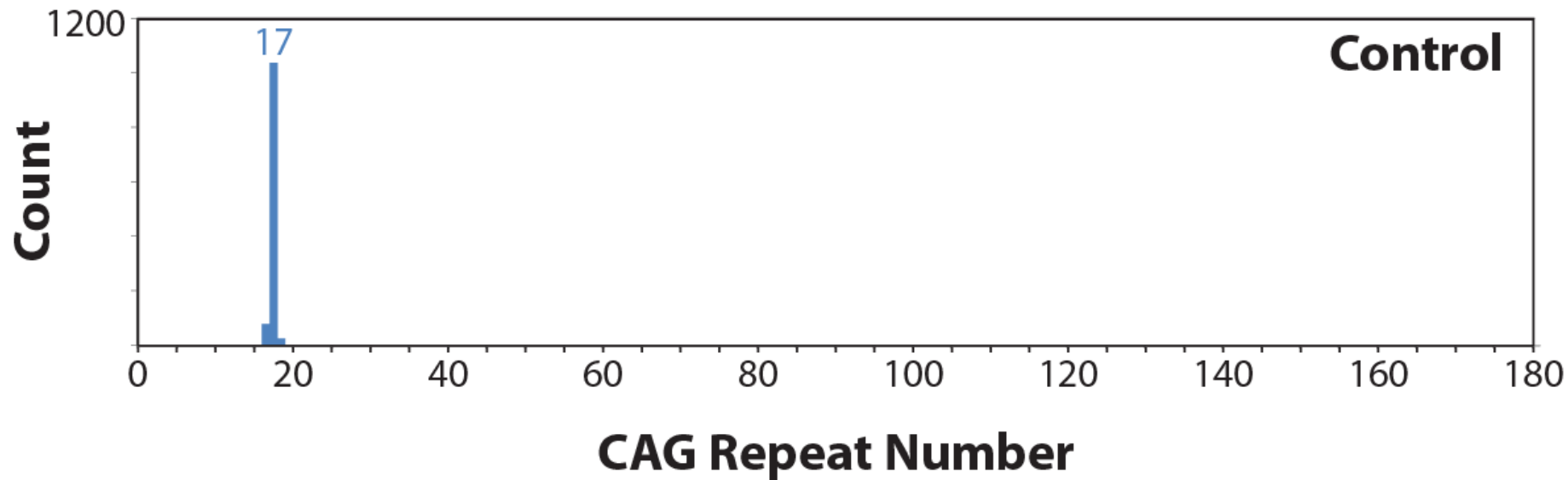
# HUNTINGTON'S DISEASE (HD)

- Autosomal dominant neurodegenerative genetic disorder
- Caused by an expansion of a CAG triplet repeat stretch in the Huntingtin (HTT) gene
  - polyglutamine tract

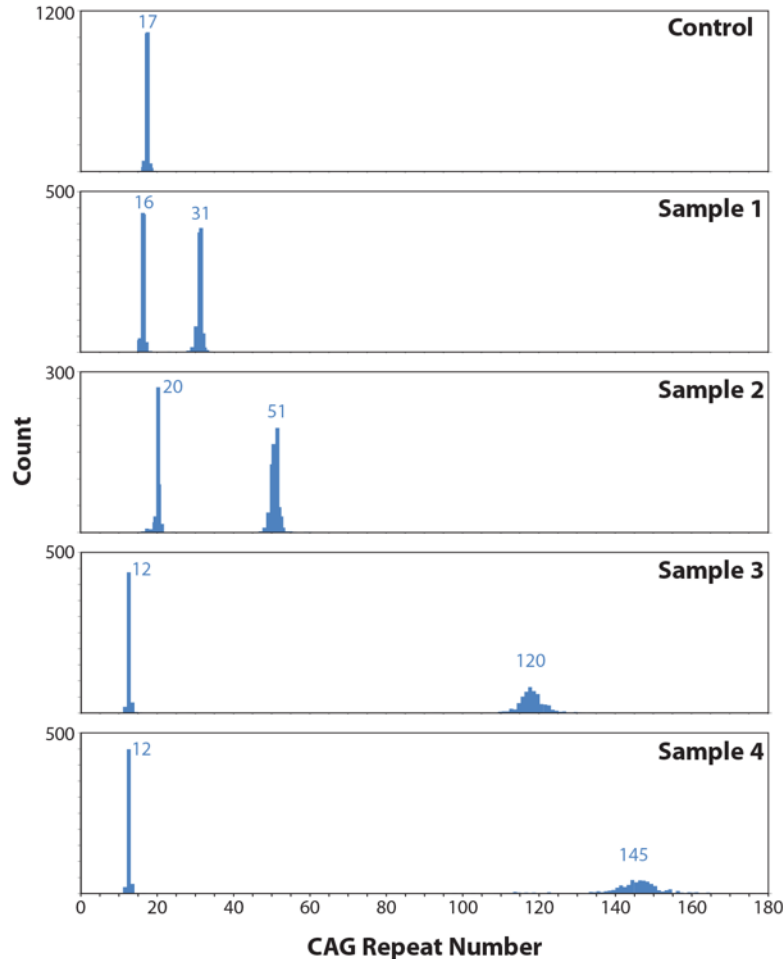




# CAG REPEAT COUNTS

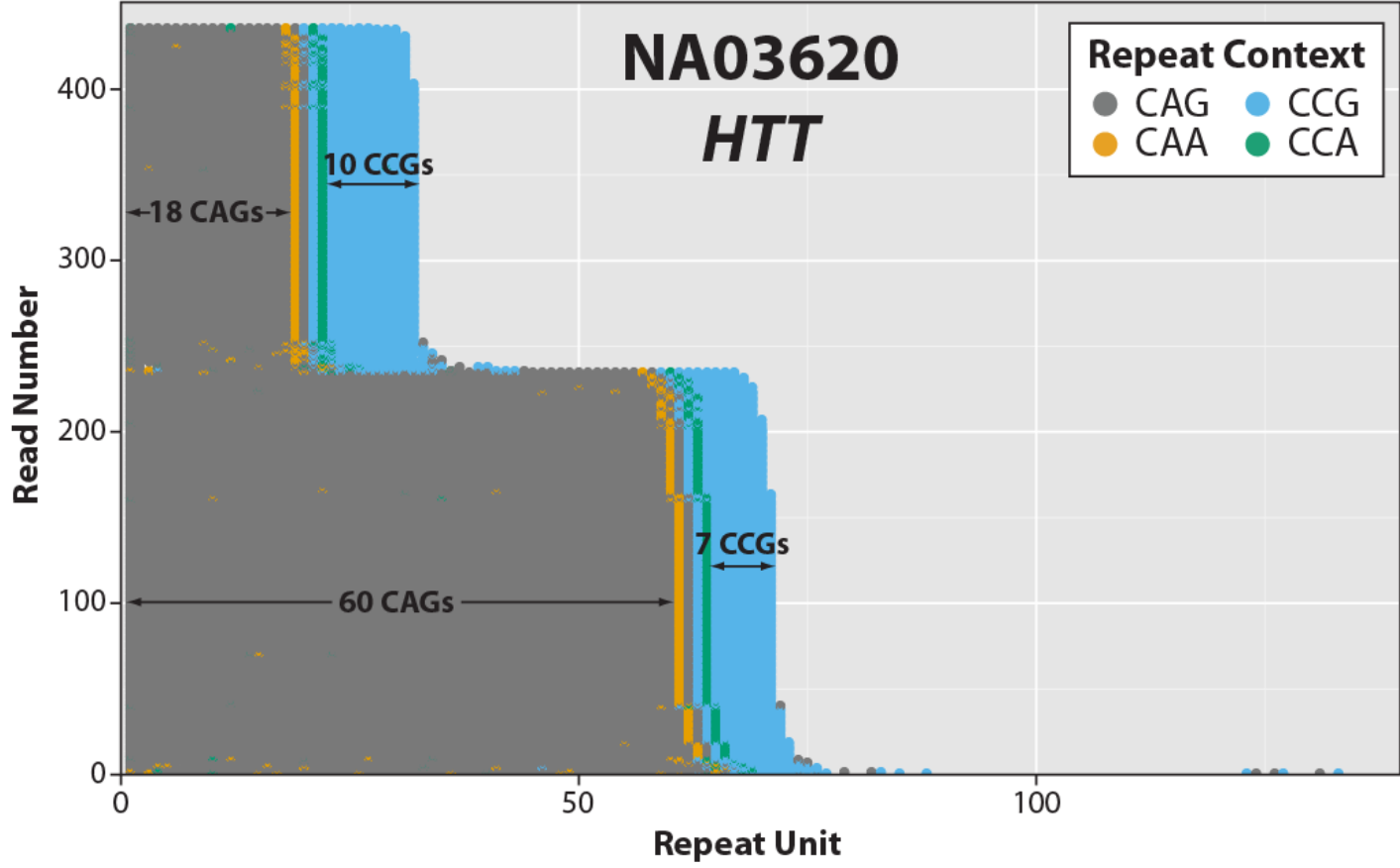


# CAG REPEAT COUNTS IN HD PATIENTS



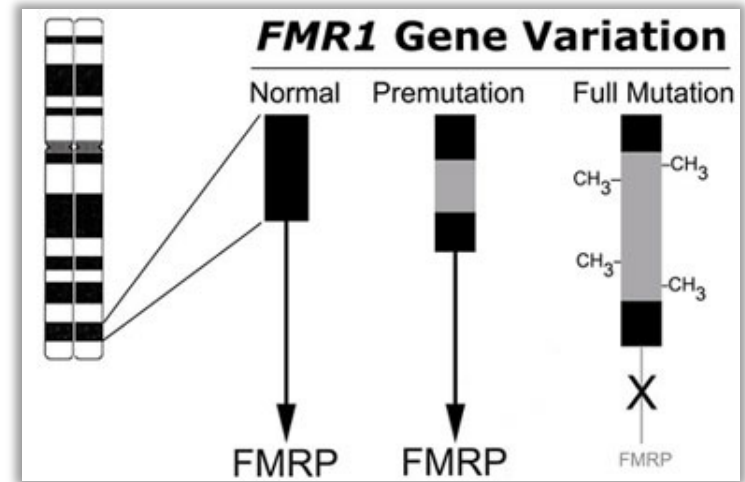
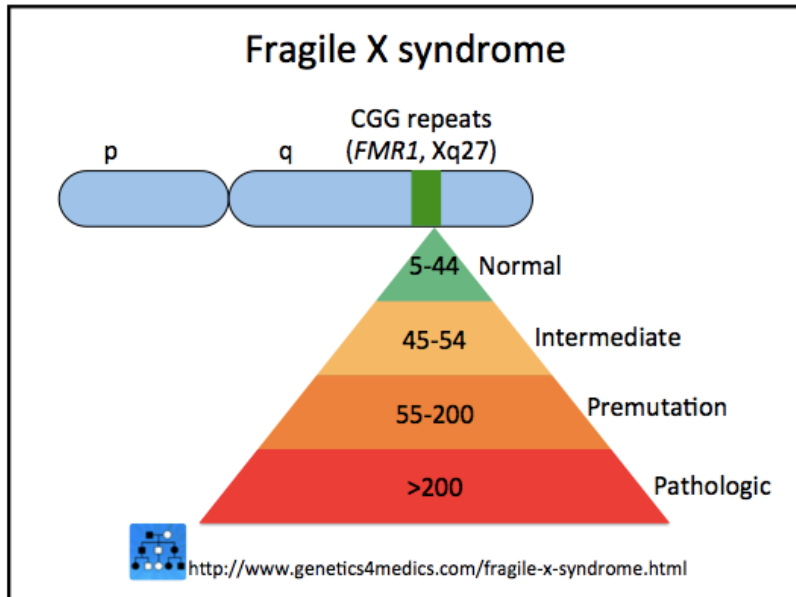
- Widening repeat number distribution at the mutated allele is biological
- Obtained roughly equal number of sequenced molecules for normal and mutated alleles

*Samples obtained from Vanessa Wheeler  
(Harvard Medical School)*



# FRAGILE X SYNDROME

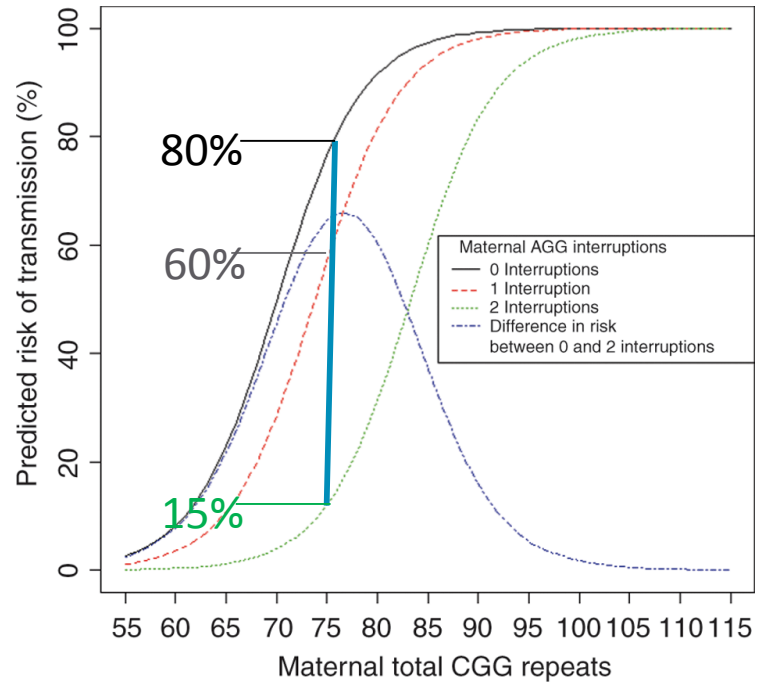
- Most common heritable form of cognitive impairment
- Caused by expansion of a CGG trinucleotide repeat in the 5' UTR of the FMR1 gene



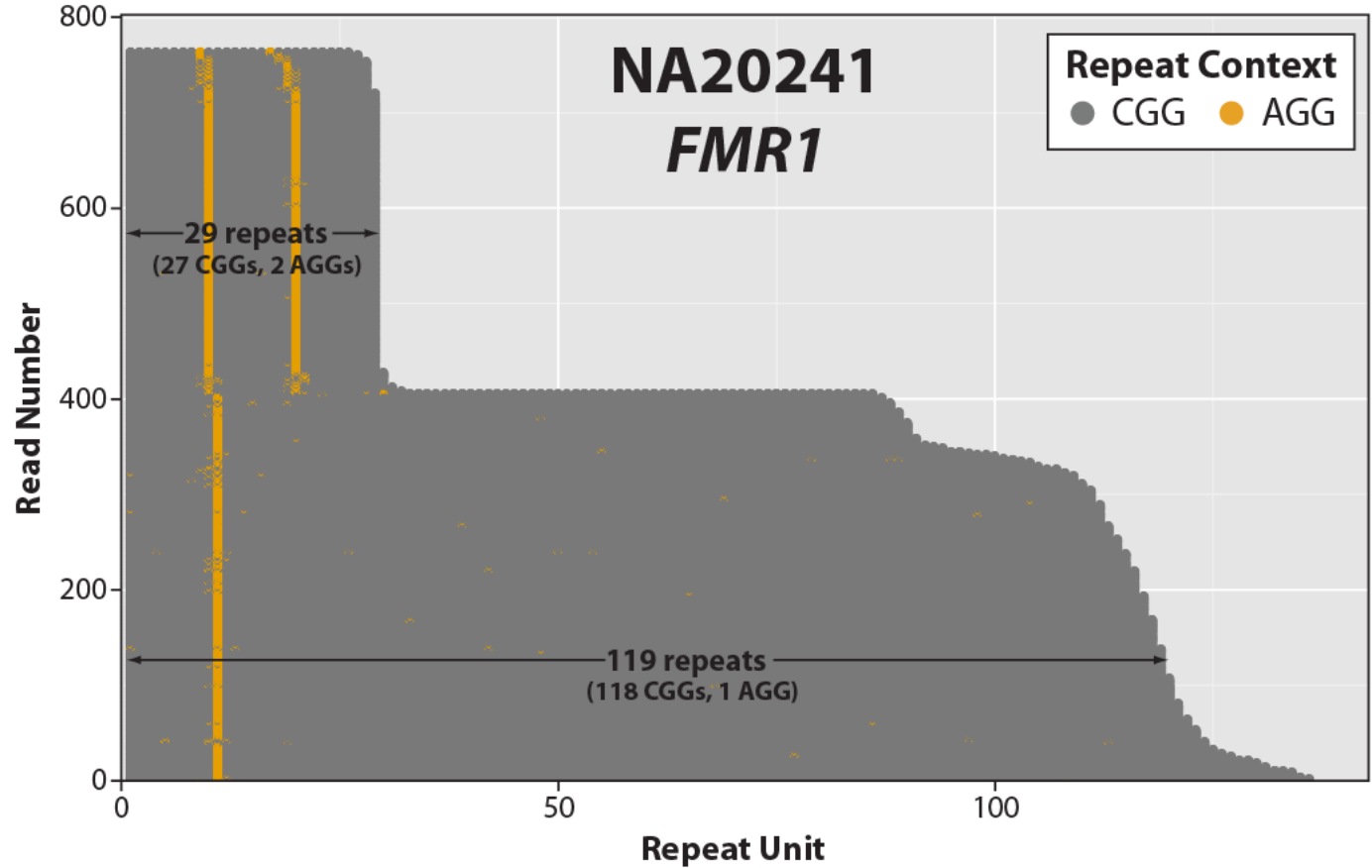
# AGG “INTERRUPTIONS” REDUCE THE CHANCES OF PRE- TO FULL MUTATION TRANSMISSION

...CGG CGG CGG CGG **AGG** CGG...

- Difference in risk is greatest near 75-80 CGG repeats
- Having full sequence information is medically relevant

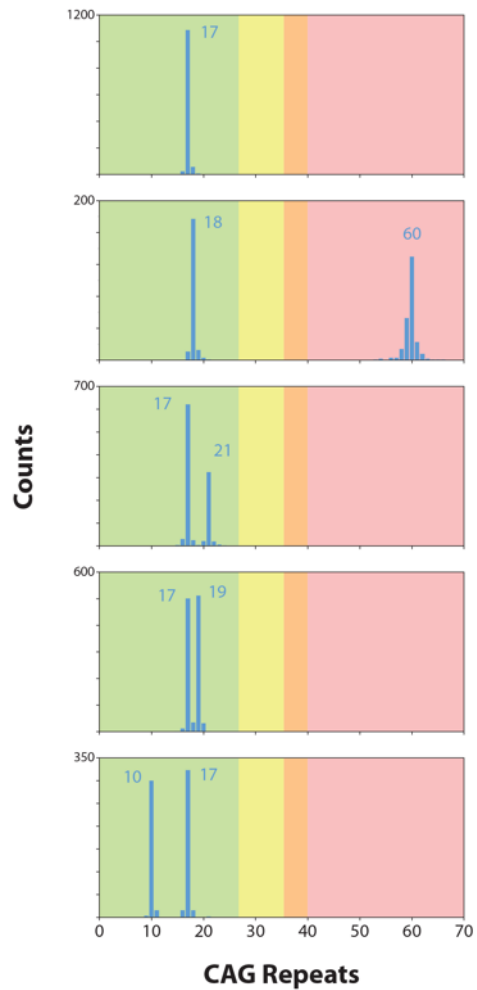


2 ...CGG CGG CGG CGG **AGG** CGG CGG CGG CGG CGG CGG CGG CGG **AGG** CGG ...  
 1 ...CGG CGG CGG CGG **AGG** CGG CGG CGG CGG CGG CGG CGG CGG CGG CGG CGG ...  
 0 ...CGG CGG CGG CGG CGG CGG CGG CGG CGG CGG CGG CGG CGG CGG CGG ...

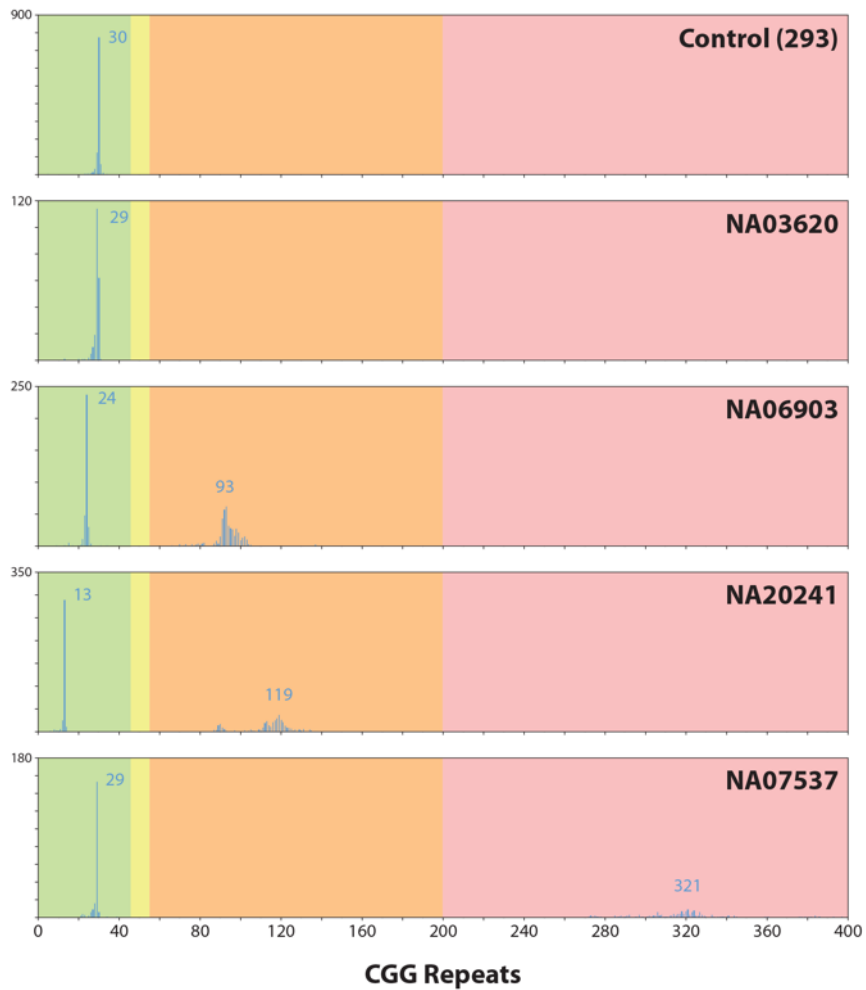




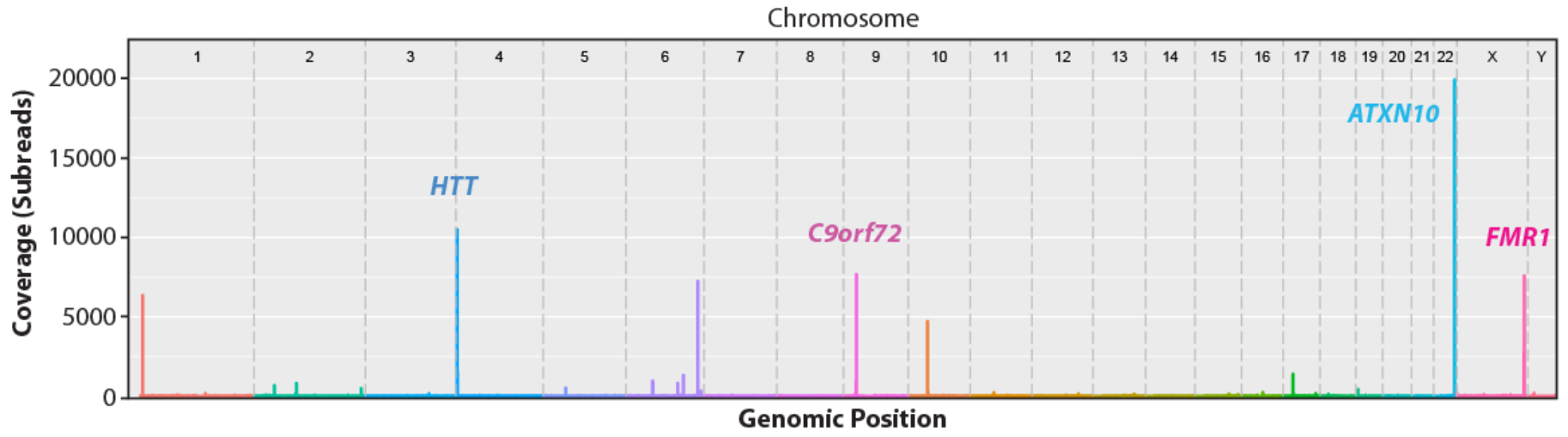
# HTT



# FMR1



# SUBREAD COVERAGE ON THE SEQUEL SYSTEM

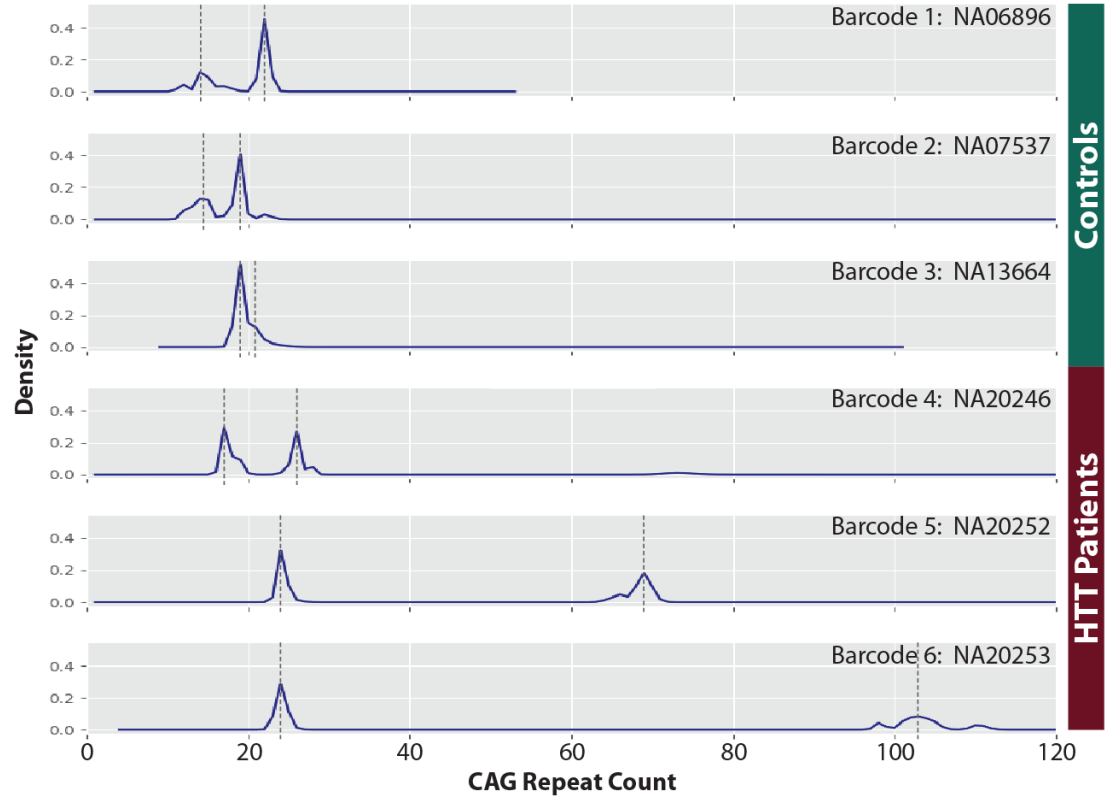


1 Sequel SMRT Cell 1M

# MULTIPLEXED SAMPLES ON THE SEQUEL SYSTEM

## HTT

CAG Repeat Counts  
from 3 Controls and  
3 HD Patients



## CONCLUSION

*Amplification-free enrichment with CRISPR/Cas9 and SMRT Sequencing achieves the base-level resolution required to understand the underlying biology of repeat expansion disorders*

- Target any hard-to-amplify genomic region regardless of sequence context
- Avoid PCR bias and PCR errors
- Accurately sequence through long repetitive and low-complexity regions
  - Count repeats and identify sequence interruptions
- Detect sample mosaicism



PACBIO®

[www.pacb.com](http://www.pacb.com)

For Research Use Only. Not for use in diagnostics procedures. © Copyright 2017 by Pacific Biosciences of California, Inc. All rights reserved. Pacific Biosciences, the Pacific Biosciences logo, PacBio, SMRT, SMRTbell, Iso-Seq, and Sequel are trademarks of Pacific Biosciences. BluePippin and SageELF are trademarks of Sage Science. NGS-go and NGSengine are trademarks of GenDx. FEMTO Pulse and Fragment Analyzer are trademarks of Advanced Analytical Technologies.

All other trademarks are the sole property of their respective owners.