

Lecture 8. Structural Variation Analysis

Michael Schatz

Feb 23, 2017

JHU 600.649: Applied Comparative Genomics



Assignment I: Due Thursday @ 11:59pm

Email PDF to: jhuappliedgenomics@gmail.com

The screenshot shows a GitHub repository page for 'appliedgenomics'. The repository has 5 stars, 8 forks, and 0 issues. The README.md file is displayed, containing instructions for Assignment 1: Genome Assembly. The assignment is due on Feb. 23, 2017, at 11:59pm. It requires coverage analysis and assembly of unassembled reads from a mysterious pathogen. Tools like Allpaths are mentioned as not working on Mac. A link to download reads and reference genome is provided.

schatzlab / appliedgenomics

Branch: master [appliedgenomics / assignments / assignment1 / README.md](#)

mschatz Update README.md 31eccf2 10 days ago

1 contributor

138 lines (96 sloc) 8.07 KB

Assignment 1: Genome Assembly

Assignment Date: Thursday, Feb. 9, 2017
Due Date: Thursday, Feb. 23, 2017 @ 11:59pm

Assignment Overview

In this assignment, you are given a set of unassembled reads from a mysterious pathogen that contains a secret message encoded someplace in the genome. The secret message will be recognizable as a novel insertion of sequence not found in the reference. Your task is to assess the quality of the reads, assemble the genome, identify, and decode the secret message. If all goes well the secret message should decode into a recognizable english text, otherwise doublecheck your coordinates and try again. As a reminder, any questions about the assignment should be posted to [Piazza](#)

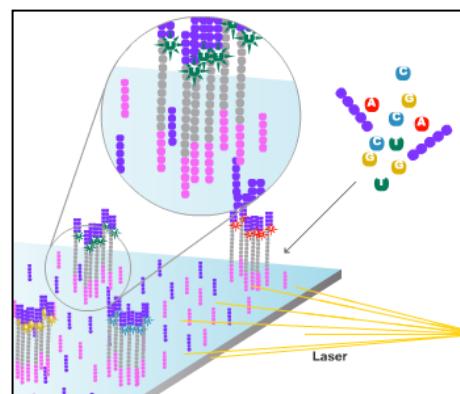
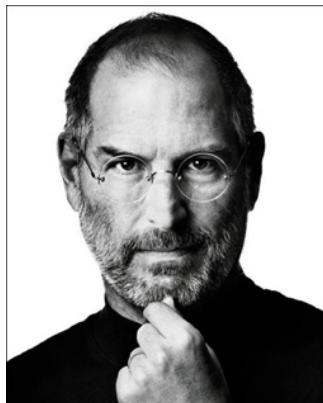
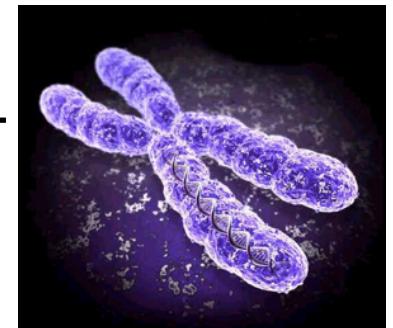
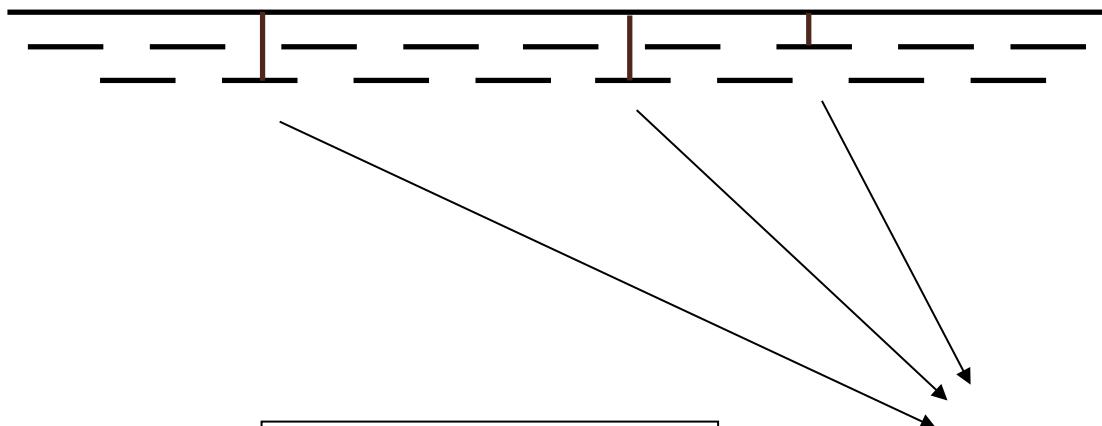
Some of the tools you will need to use only run in a linux environment. Allpaths, for example, will *not* work under Mac, even though it will compile. If you do not have access to a linux machine, download and install a virtual machine following the directions here: <https://github.com/schatzlab/appliedgenomics/blob/master/assignments/virtualbox.md>

Question 1. Coverage Analysis [10 pts]

Download the reads and reference genome from:
<https://github.com/schatzlab/appliedgenomics/raw/master/assignments/assignment1/asm.tgz>

Personal Genomics

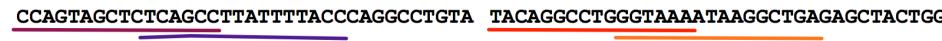
How does your genome compare to the reference?



Heart Disease
Cancer
Creates magical
technology

Read Mapping Overview

1. Split read into segments

Read

Policy: extract 16 nt seed every 10 nt

Seeds

+ , 0: CCAGTAGCTCTCAGCC	- , 0: TACAGGCCTGGGTAAA
+ , 10: TCAGCCTTATTTACC	- , 10: GGTAAAAATAAGGCTGA
+ , 20: TTTACCCAGGCCTGTA	- , 20: GGCTGAGAGCTACTGG

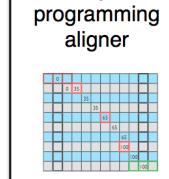
2. Lookup each segment and prioritize

Seeds

+ , 0: CCAGTAGCTCTCAGCC	→	Ungapped alignment with FM Index
+ , 10: TCAGCCTTATTTACC	→	a a c a a c a c a c a c \$ a c g \$ a c a I c - - - a g \$ a c a c
+ , 20: TTTACCCAGGCCTGTA	→	Seed alignments (as B ranges)
- , 0: TACAGGCCTGGGTAAA	→	{ [211, 212], [212, 214] } { [653, 654], [651, 653] }
- , 10: GGTAAAAATAAGGCTGA	→	{ [684, 685] }
- , 20: GGCTGAGAGCTACTGG	→	{ }
	→	{ }
	→	{ [624, 625] }

3. Evaluate end-to-end match

Extension candidates

SA:684, chr12:1955	→	SIMD dynamic programming aligner
SA:624, chr2:462	→	
SA:211: chr4:762	→	
SA:213: chr12:1935	→	
SA:652: chr12:1945	→	

SIMD dynamic programming aligner

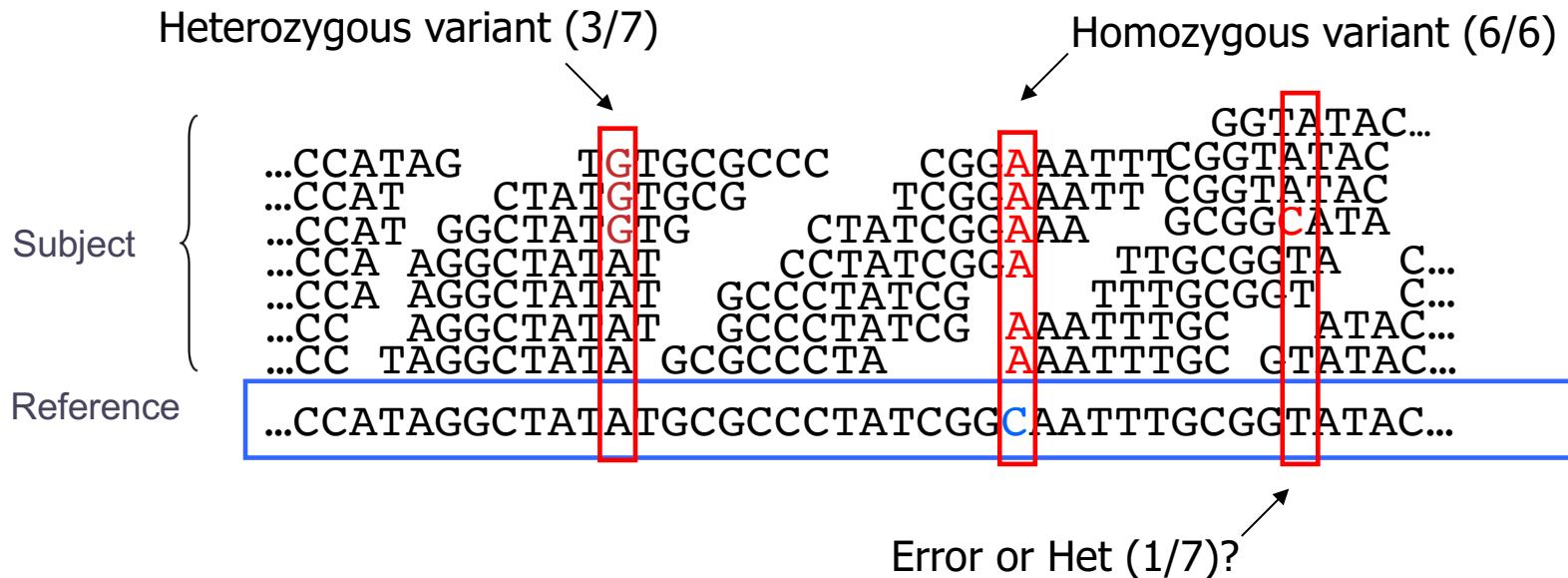
SAM alignments

r1	0	chr12	1936	0
	36M	*	0	0
		CCAGTAGCTCTCAGCCTTATTTACCCAGGCCTGTA		
		II		
		AS:i:0 XS:i:-2 XN:i:0		
		XM:i:0 XO:i:0 XG:i:0		
		NM:i:0 MD:Z:36 YT:Z:UU		
		YM:i:0		
		...		

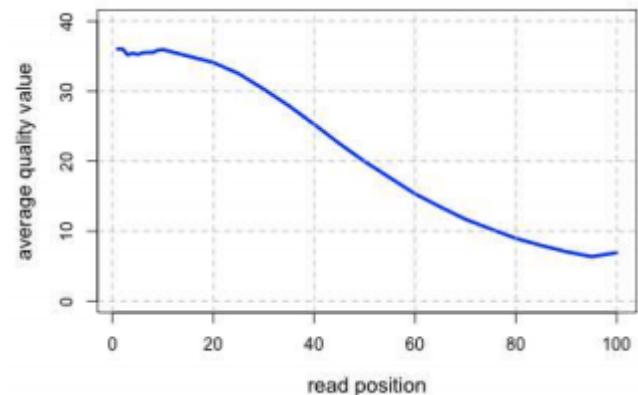
Fast gapped-read alignment with Bowtie 2

Langmead & Salzberg (2012) Nature Methods. doi:10.1038/nmeth.1923

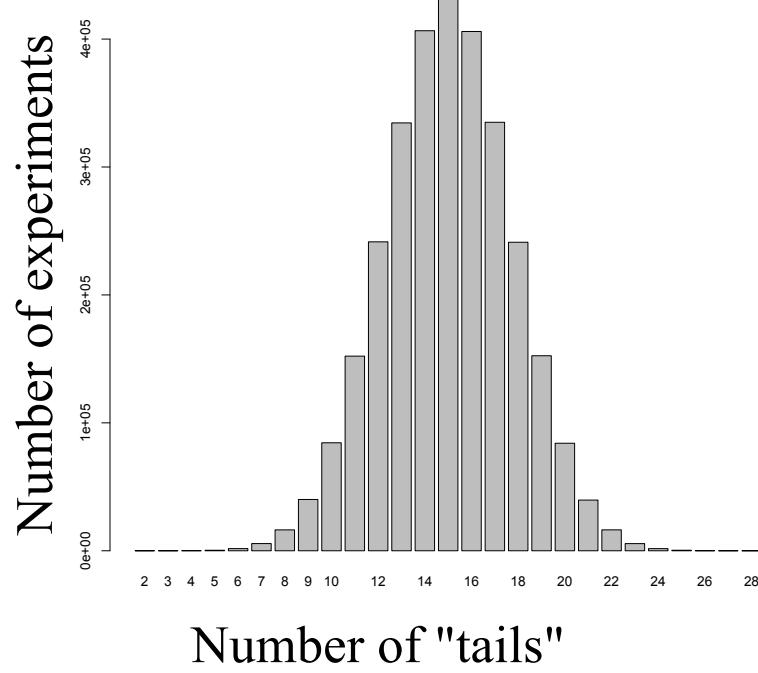
Genotyping Theory



- If there were no sequencing errors, identifying SNPs would be very easy: any time a read disagrees with the reference, it must be a variant!
- Sequencing instruments make mistakes
 - Quality of read decreases over the read length
- A single read differing from the reference is probably just an error, but it becomes more likely to be real as we see it multiple times



So, with 30 tosses (reads), we are much more likely to see an even mix of alternate and reference alleles at a heterozygous locus in a genome

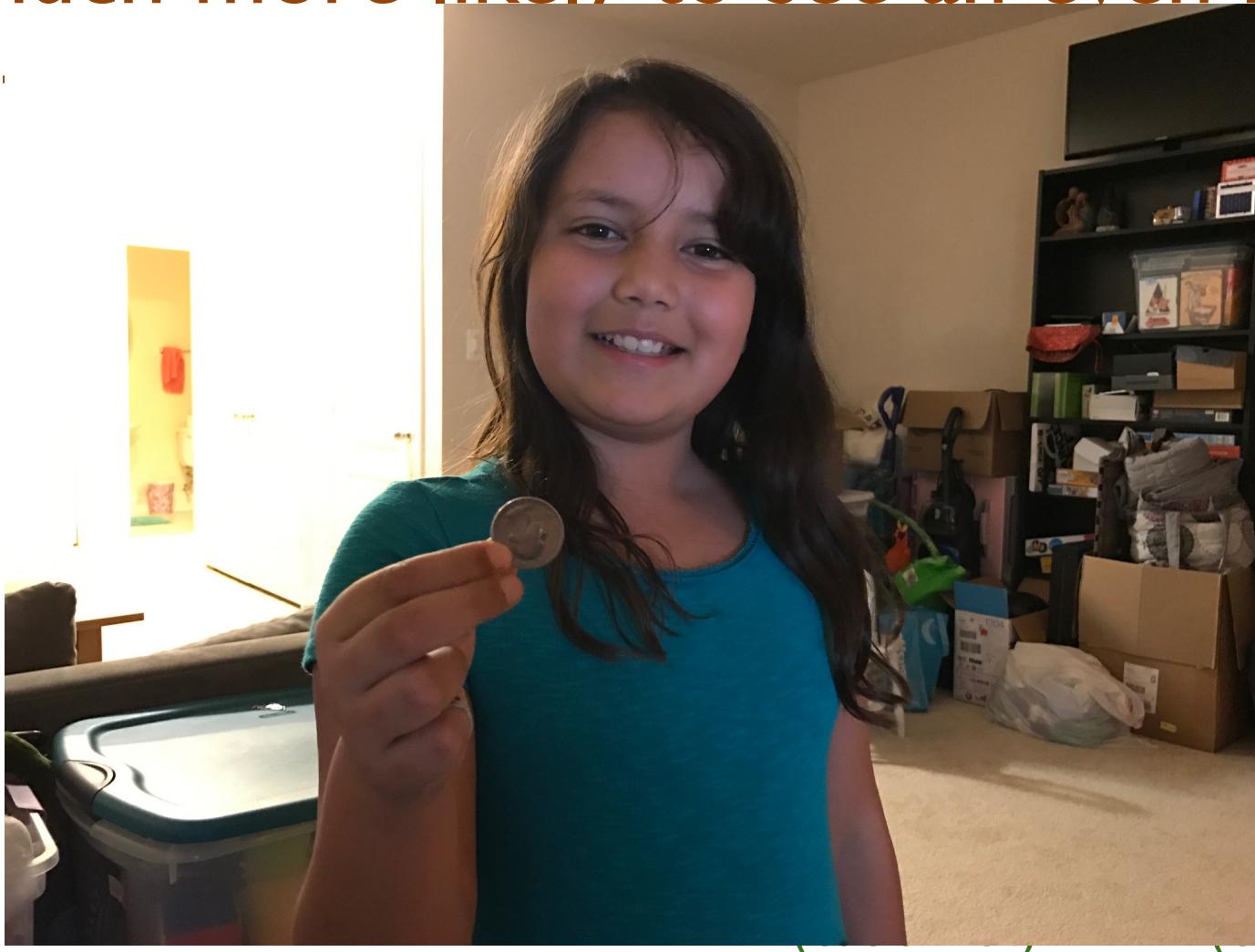


This is why at least a "30X" (30 fold sequence coverage) genome is recommended: it confers sufficient power to distinguish heterozygous alleles and from mere sequencing errors

$P(3/30 \text{ het}) <?> P(3/30 \text{ err})$

So, with 30 tosses (reads), we are much more likely to see an even mix of heads/tails.

Number of experiments



a "30X" coverage) tended: it power to spurious mere errors /30 err)

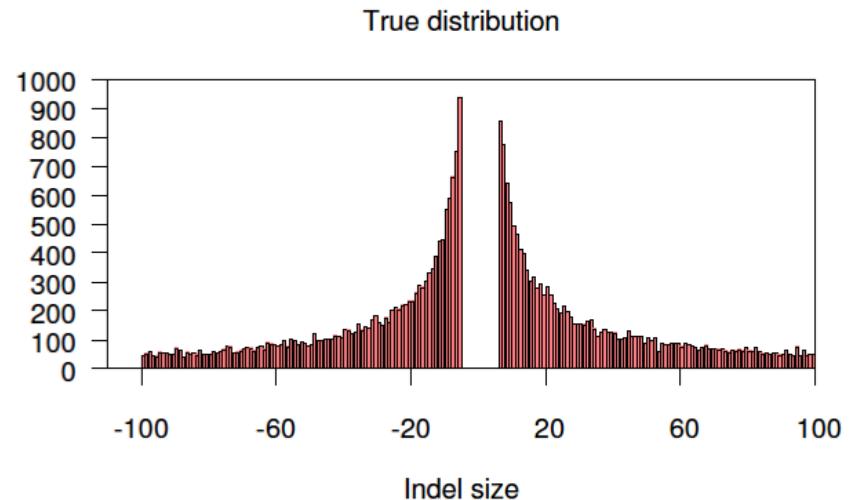
Variation Detection Complexity

SNPs + Short Indels

High precision and sensitivity

.. TTTAGAATAG-CGAGTGC ...

| | | | | | | | | | | | | |
AGAATAG**G**CGAG

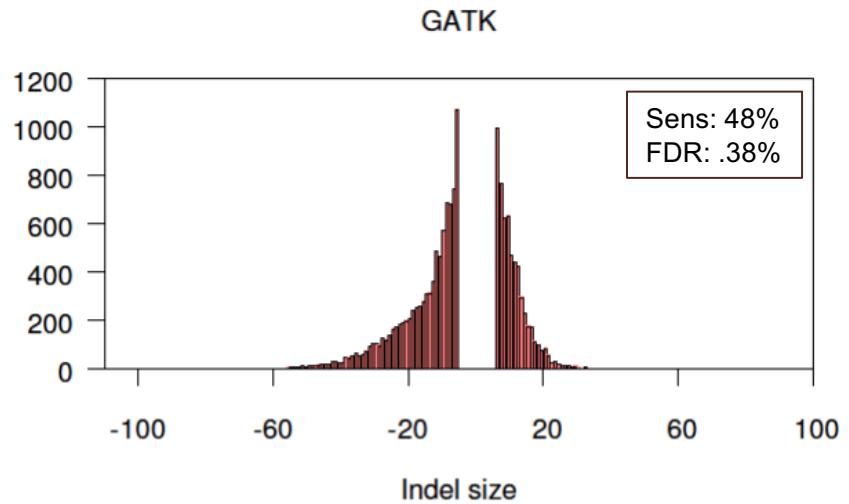


“Long” Indels (>5bp)

Reduced precision and sensitivity

.. TTTAG-----AGTGC ...

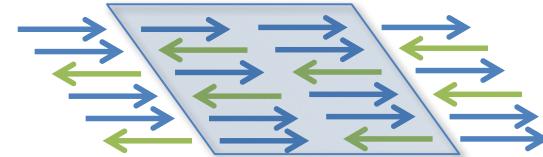
| | | | | | | | | | | | | |
TTTAG**AATAGGC** | | | | | |
ATAGGCGAGTGC



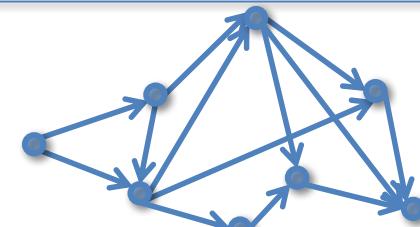
Analysis confounded by sequencing errors, localized repeats, allele biases, and mismapped reads

Scalpel Algorithm

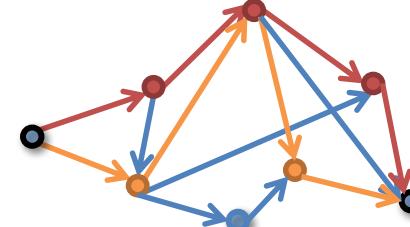
Extract reads mapping within the exon including (1) well-mapped reads, (2) soft-clipped reads, and (3) anchored pairs



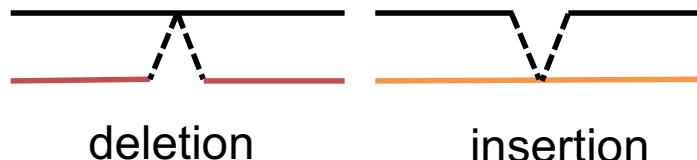
Decompose reads into overlapping k -mers and construct de Bruijn graph from the reads



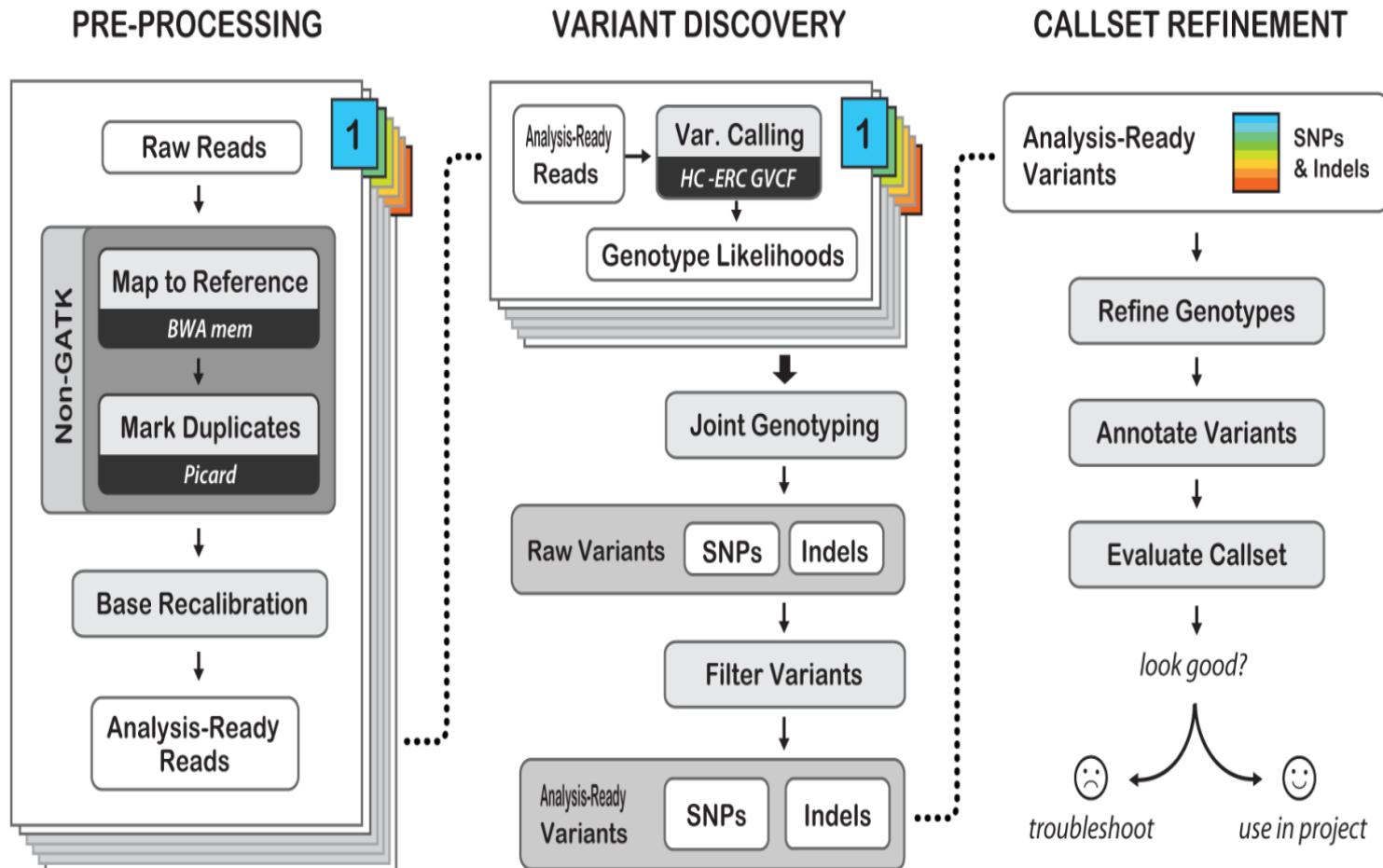
Find end-to-end haplotype paths spanning the region



Align assembled sequences to reference to detect mutations

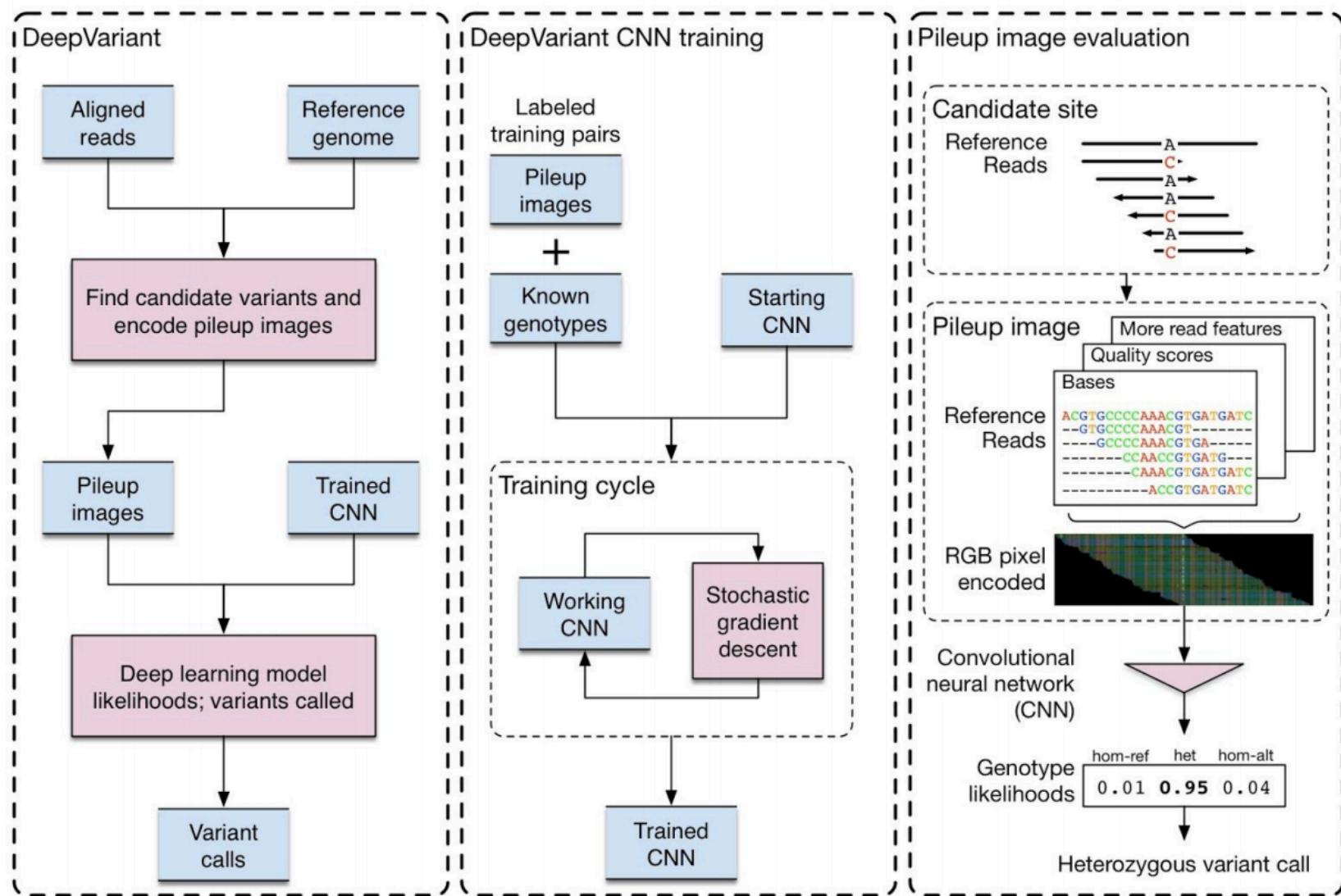


GATK workflow



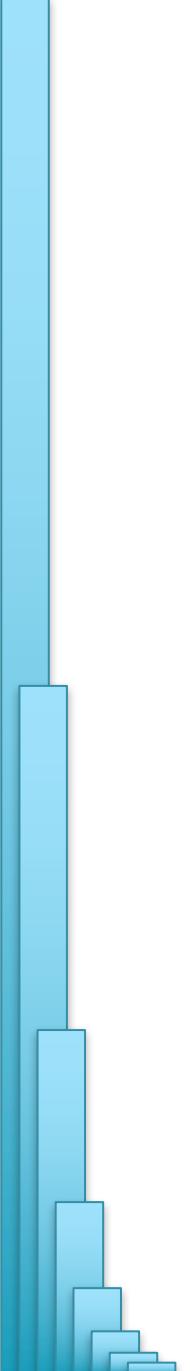
Best Practices for Germline SNPs and Indels in Whole Genomes and Exomes - June 2016

Deep Variant



Creating a universal SNP and small indel variant caller with deep neural networks

Poplin et al. (2016) bioRxiv. doi: <https://doi.org/10.1101/092890>

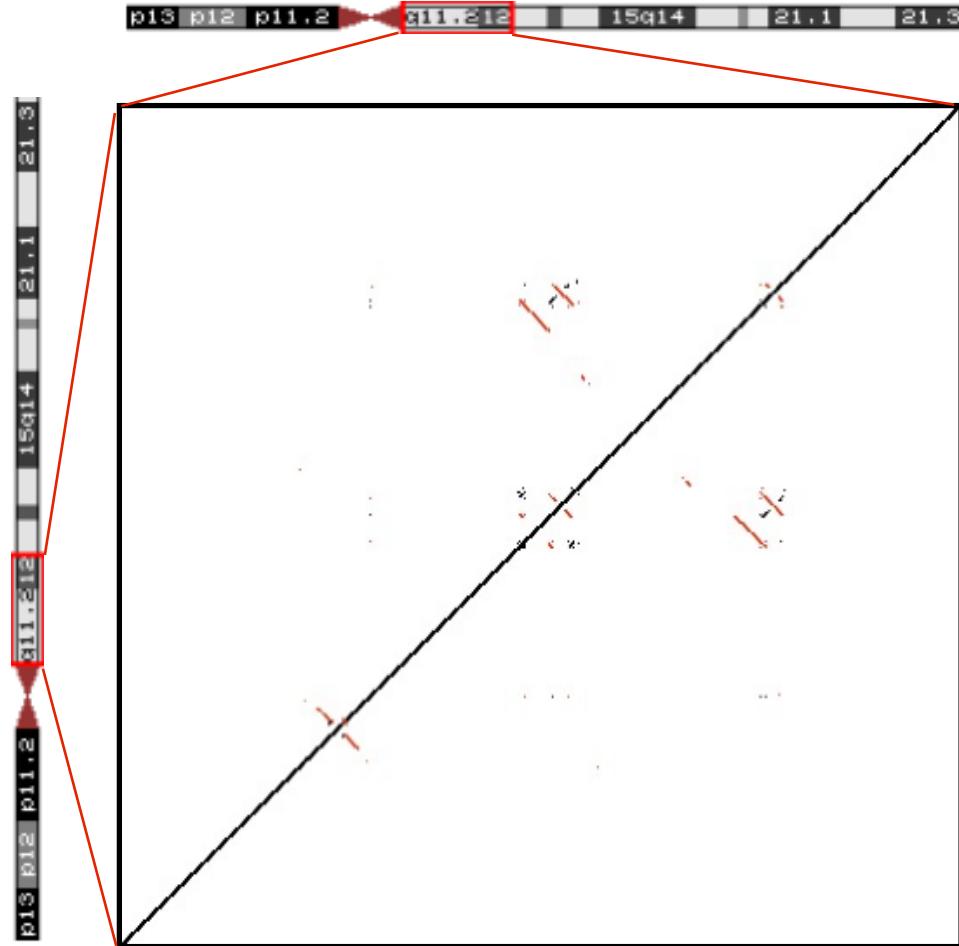


Part II: Structural Variations

Early 2000s dogma: SNPs account for most human genetic variation



Segmental duplications (a.k.a. Low copy repeats)

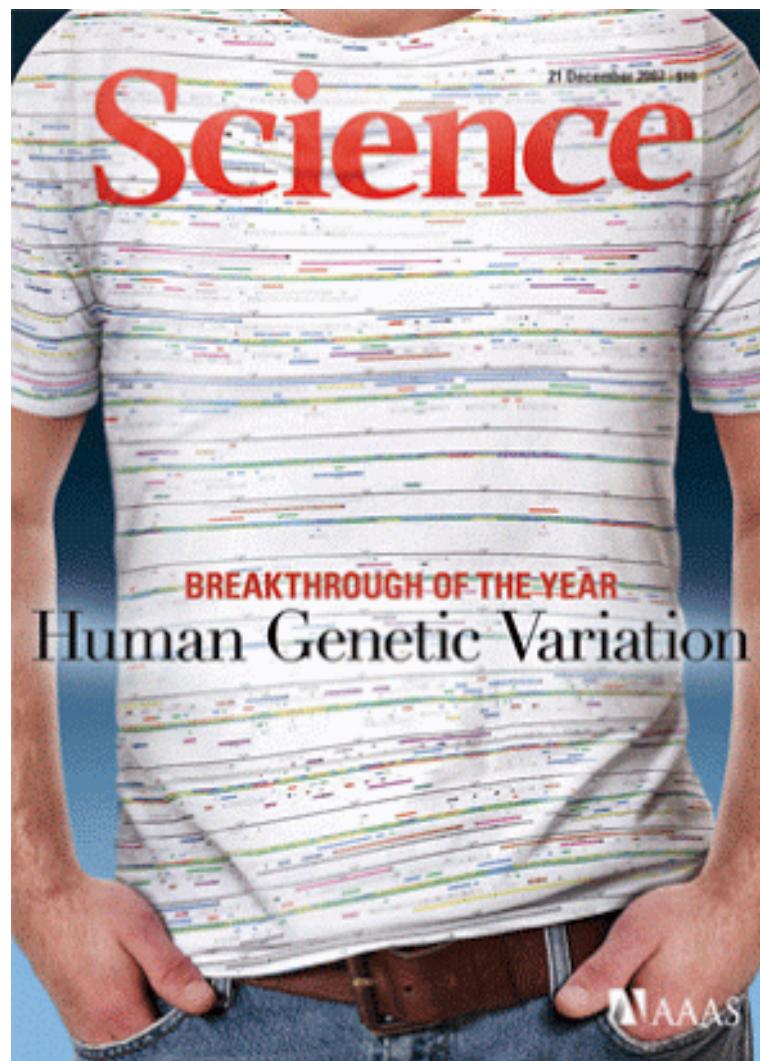


Self Dotplot:
10 megabases of Chr 15
(dot = 1 kb exact match)

~5% of the human genome is duplicated!

Bailey et al, 2002

2007 Science magazine breakthrough of the year? Why?



Discovery of abundant copy-number variation

Science, July 2004

Large-Scale Copy Number Polymorphism in the Human Genome

Jonathan Sebat,¹ B. Lakshmi,¹ Jennifer Troge,¹ Joan Alexander,¹ Janet Young,² Pär Lundin,³ Susanne Mänér,³ Hillary Massa,² Megan Walker,² Maoyen Chi,¹ Nicholas Navin,¹ Robert Lucito,¹ John Healy,¹ James Hicks,¹ Kenny Ye,⁴ Andrew Reiner,¹ T. Conrad Gilliam,⁵ Barbara Trask,² Nick Patterson,⁶ Anders Zetterberg,³ Michael Wigler^{1*}

76 CNVs in 20 individuals

70 genes

Nature Genetics, Aug. 2004

Detection of large-scale variation in the human genome

A John Iafrate^{1,2}, Lars Feuk³, Miguel N Rivera^{1,2}, Marc L Listewnik¹, Patricia K Donahoe^{2,4}, Ying Qi³, Stephen W Scherer^{3,5} & Charles Lee^{1,2,5}

255 CNVs in 55 individuals

127 genes

- 331 CNVs, only 11 in common
- Half observed in only 1 individual
- Impact "plenty" of genes
- Correlated with segmental duplications in the reference genome

Why is structural variation relevant / important?

- They are common and affect a large fraction of the genome
 - In total, SVs impact more base pairs than all single-nucleotide differences.
- They are a major driver of genome evolution
 - Speciation can be driven by rapid changes in genome architecture
 - Genome instability and aneuploidy: hallmarks of solid tumor genomes

Why is structural variation relevant / important?

- Genetic basis of traits
 - Gene dosage effects.
 - Neuropsychiatric disease (e.g., autism, schizophrenia)
 - Spontaneous SVs implicated in so-called “genomic” and developmental disorders
 - Somatic genome instability; age-dependent disease

SV and human disease phenotypes

Mendelian (X-linked)

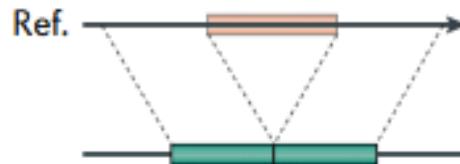
Hemophilia A	306700	<i>F8</i>	inv/del
Hunter syndrome	309900	<i>IDS</i>	del/inv
Ichthyosis	308100	<i>STS</i>	del
Mental retardation	300706	<i>HUWE1</i>	dup
Pelizaeus-Merzbacher disease	312080	<i>PLP1</i>	del/dup/tri
Progressive neurological symptoms (MR+SZ)	300260	<i>MECP2</i>	dup
Red-green color blindness	303800	opsin genes	del

Complex traits

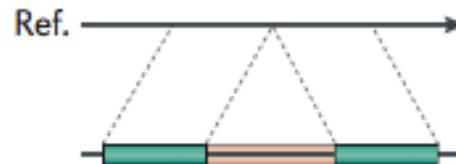
Alzheimer disease	104300	<i>APP</i>	dup
Autism	612200	3q24	inherited homozygous del
	611913	16p11.2	del/dup
Crohn disease	266600	<i>HBD-2</i>	copy number loss
	612278	<i>IRGM</i>	del
HIV susceptibility	609423	<i>CCL3L1</i>	copy number loss
Mental retardation	612001	15q13.3	del
	610443	17q21.31	del
	300534	Xp11.22	dup
Pancreatitis	167800	<i>PRSS1</i>	tri
Parkinson disease	168600	<i>SNCA</i>	dup/tri
Psoriasis	177900	<i>DEFB</i>	copy number gain
Schizophrenia	612474	1q21.1	del
	181500	15q11.2	del
	612001	15q13.3	del
Systemic lupus erythematosus	152700	<i>FCGR3B</i>	copy number loss
	120810	<i>C4</i>	copy number loss

Structural Variations

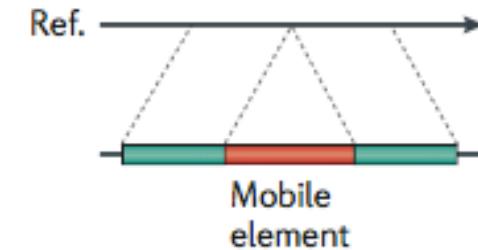
Deletion



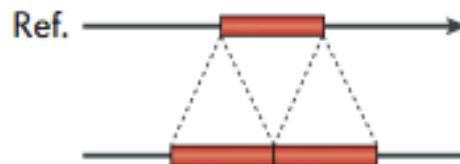
Novel sequence insertion



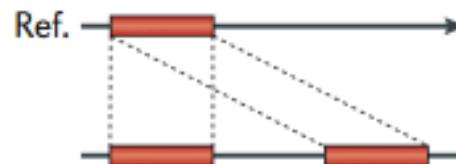
Mobile-element insertion



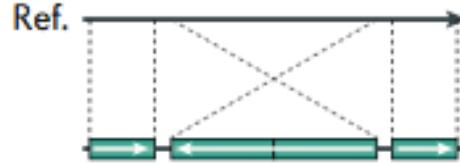
Tandem duplication



Interspersed duplication



Inversion



Translocation



**Any mutation >50bp
Profound impact on
genome structure
and function**

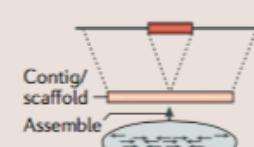
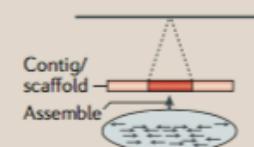
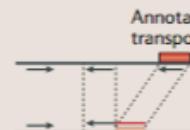
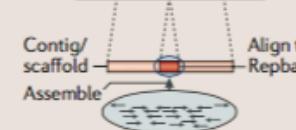
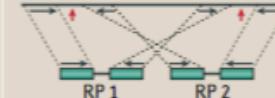
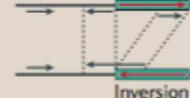
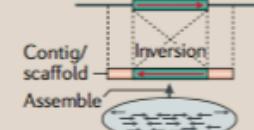
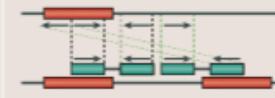
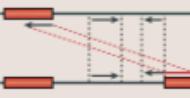
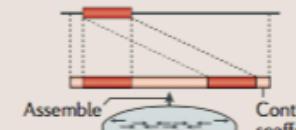
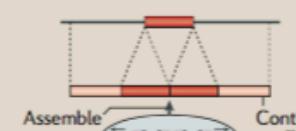
Genome structural variation discovery and genotyping

Alkan, C, Coe, BP, Eichler, EE (2011) *Nature Reviews Genetics*. May;12(5):363-76. doi: 10.1038/nrg2958.

Structural Variation Sequence Signatures

SV classes	Read pair	Read depth	Split read	Assembly
Deletion				
Novel sequence insertion		Not applicable		
Mobile-element insertion		Not applicable		
Inversion		Not applicable		
Interspersed duplication				
Tandem duplication				

Structural Variation Sequence Signatures

SV classes	Read pair	Read depth	Split read	Assembly
Deletion				
Novel sequence insertion		Not applicable		
Mobile-element insertion		Not applicable		
Inversion		Not applicable		
Interspersed duplication				
Tandem duplication				

Paired-end and Mate-pairs

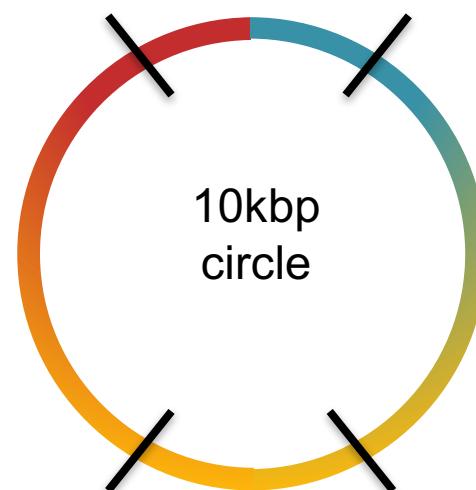
Paired-end sequencing

- Read one end of the molecule, flip, and read the other end
- Generate pair of reads separated by up to 500bp with inward orientation



Mate-pair sequencing

- Circularize long molecules (1-10kbp), shear into fragments, & sequence
- Mate failures create short paired-end reads



2x100 @ ~10kbp (outies)



2x100 @ 300bp (innies)



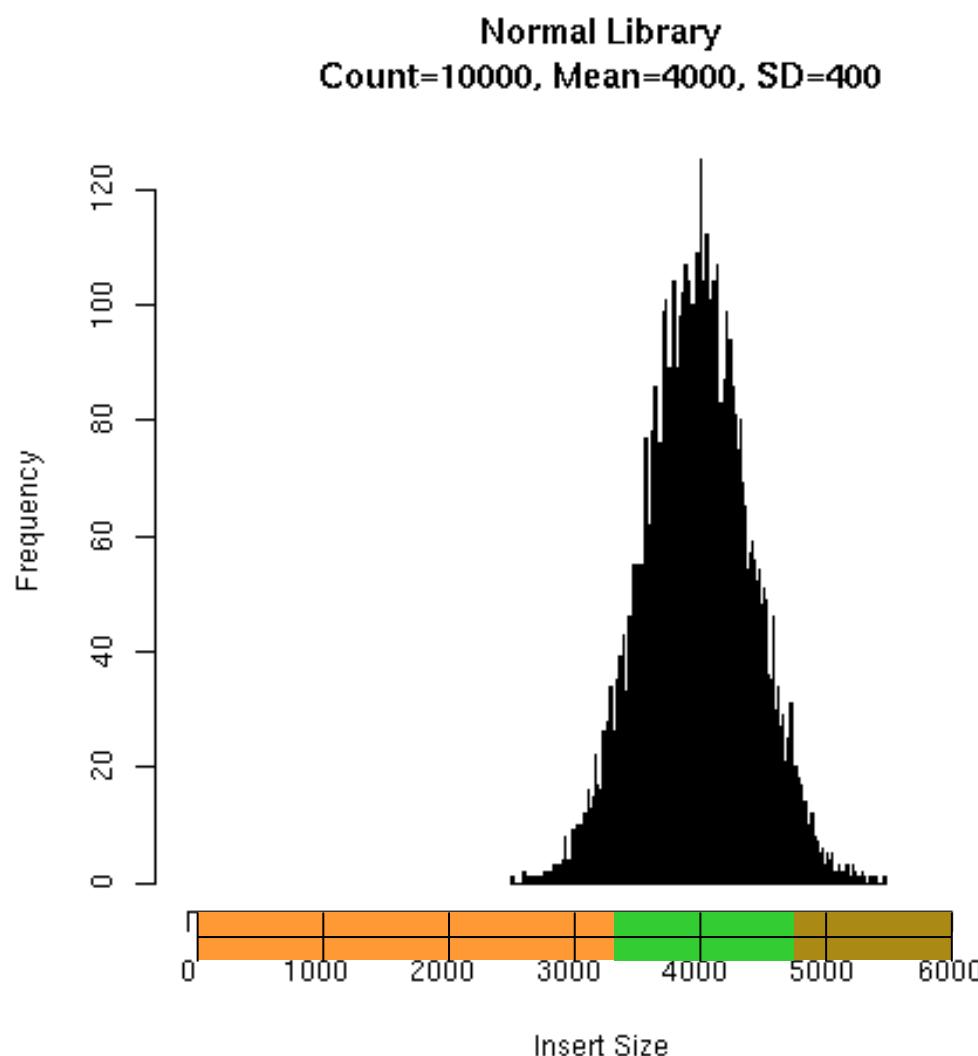
C/E Statistic

- The presence of individual compressed or expanded mates is rare but expected.
- Do the inserts spanning a given position differ from the rest of the library?
 - Flag large differences as potential structural variation / misassemblies
 - Even if each individual mate is “happy”
- Compute the statistic at all positions
 - $(\text{Local Mean} - \text{Global Mean}) / \text{Scaling Factor}$
- Introduced by Jim Yorke’s group at UMD

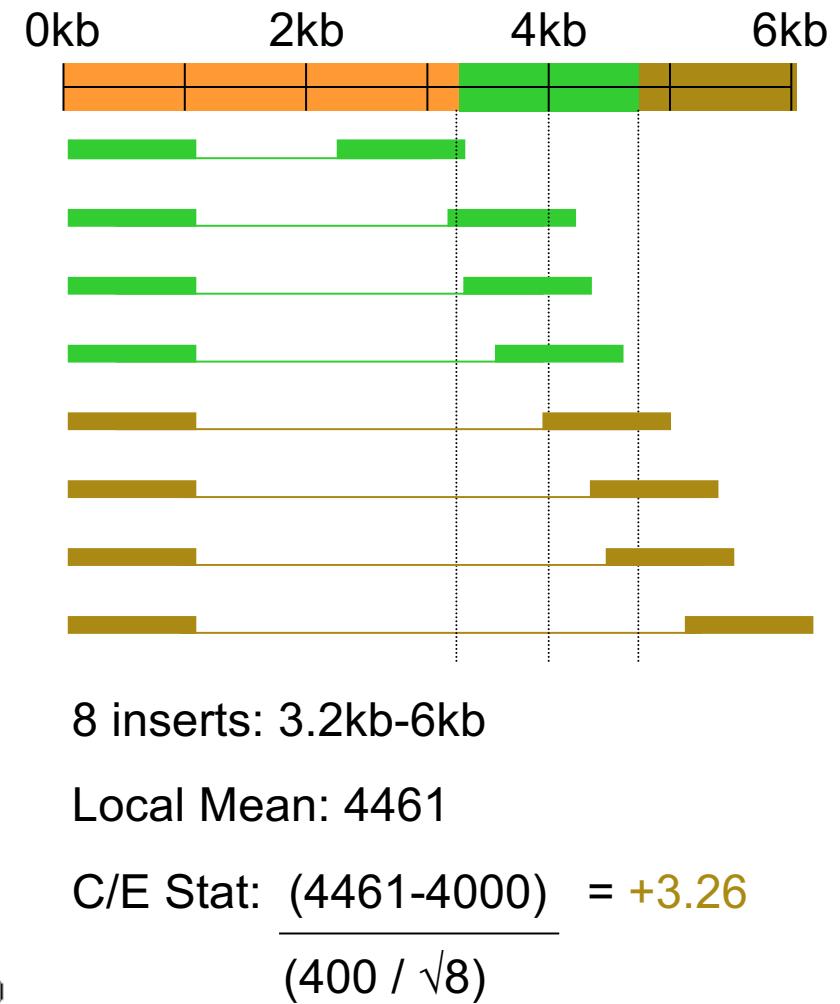
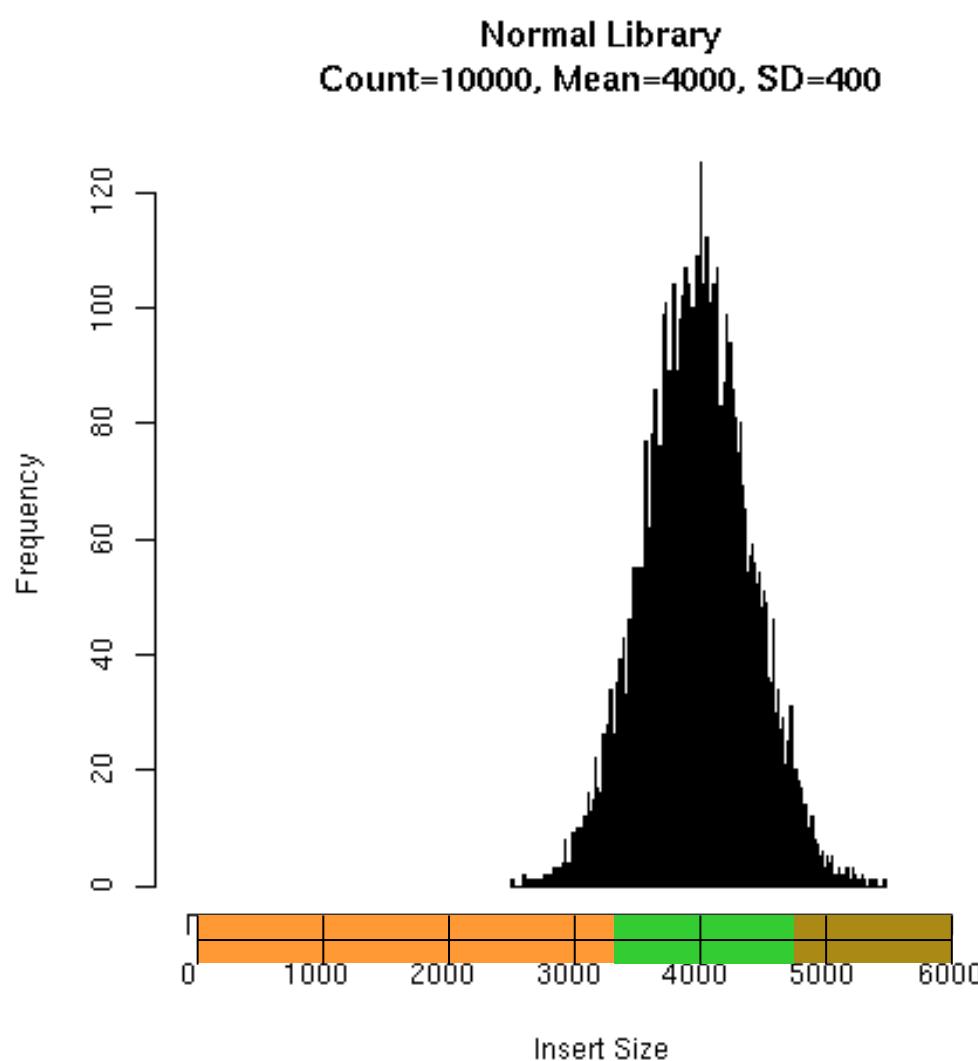
Assembly reconciliation

Zimin, Smith, Sutton, Yorke (2008) Bioinformatics. doi 10.1093/bioinformatics/btm542

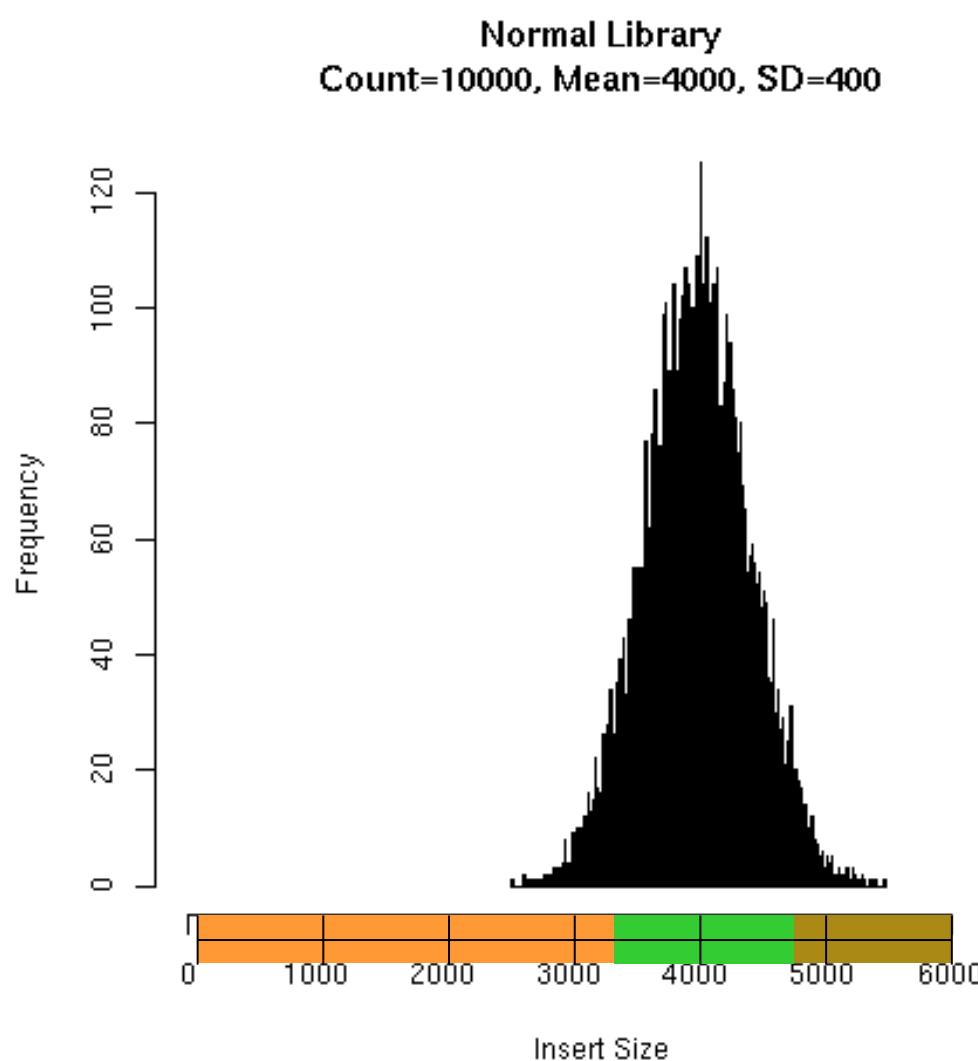
Sampling the Genome



C/E-Statistic: Expansion



C/E-Statistic: Compression



8 inserts: 3.2 kb-4.8kb

Local Mean: 3488

$$\text{C/E Stat: } \frac{(3488 - 4000)}{(400 / \sqrt{8}}) = -3.62$$

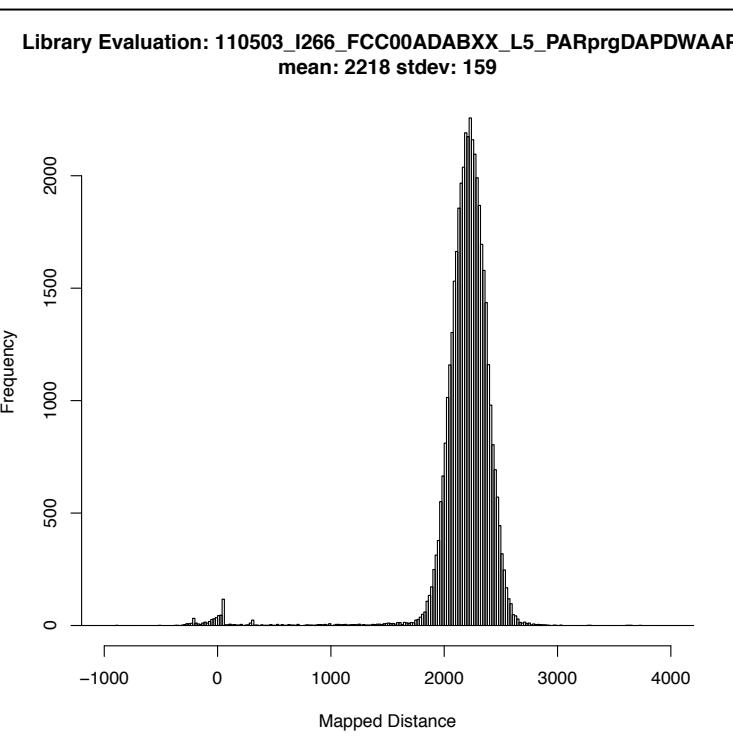
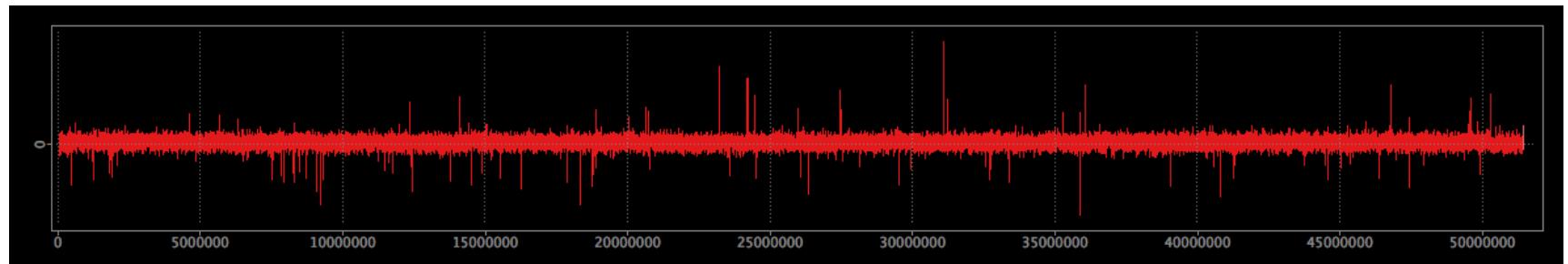
C/E Stat ≤ -3.0 indicates
Compression

Parrot Metassembly

CE statistic (projected) across 51.1 Mbp scaffold

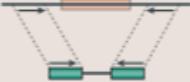
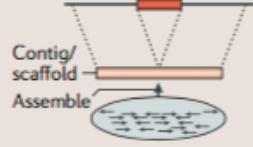
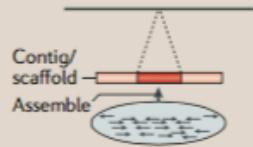
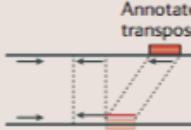
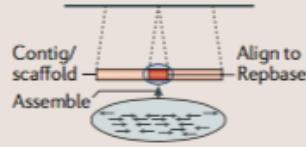
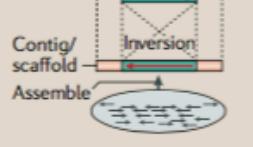
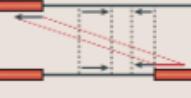
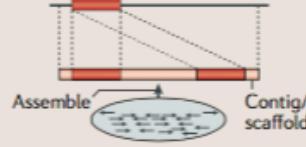
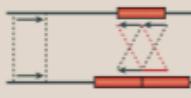
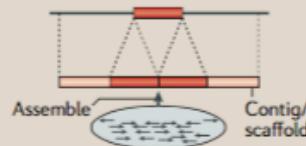


6C

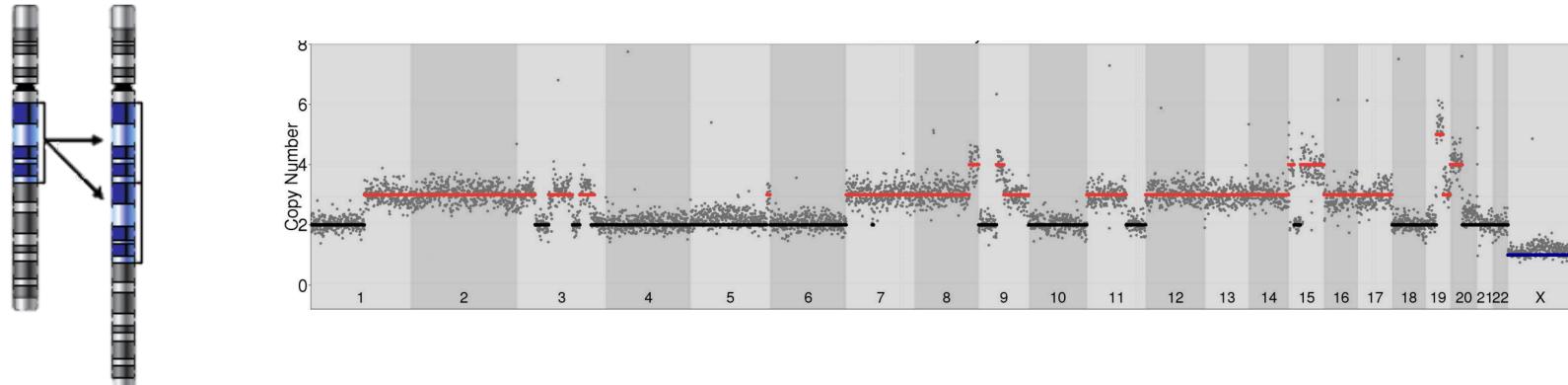
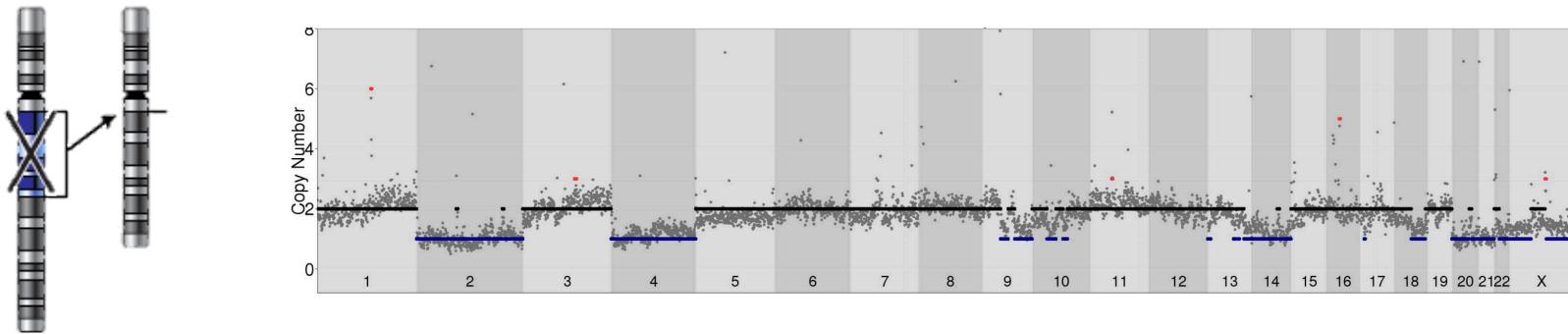
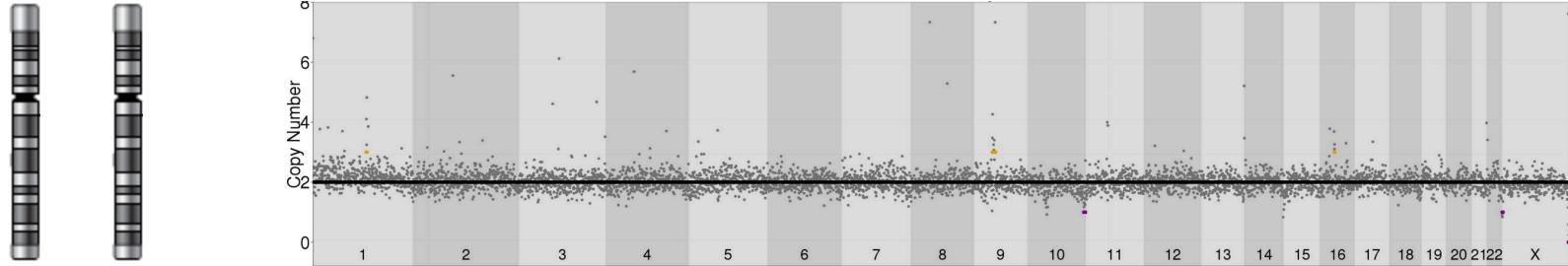


- Re-map 2kbp mates to each draft assembly, compute CE statistic at every position
- Extreme CE values are likely to be mis-assemblies
 - Can also look at coverage, mis-oriented mates, and other forensics features
 - Approximately 1.4 major events per Mbp

Structural Variation Sequence Signatures

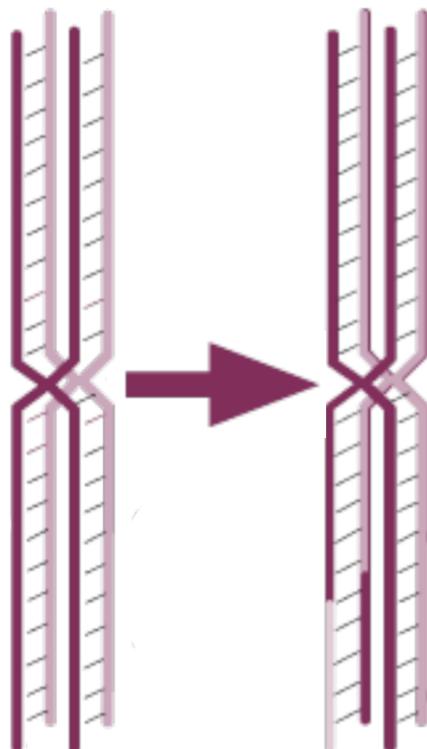
SV classes	Read pair	Read depth	Split read	Assembly
Deletion				 Contig/ scaffold Assemble
Novel sequence insertion		Not applicable		 Contig/ scaffold Assemble
Mobile-element insertion		Not applicable		 Contig/ scaffold Assemble Align to Repbase
Inversion		Not applicable		 Contig/ scaffold Assemble Inversion
Interspersed duplication				 Assemble Contig/ scaffold
Tandem duplication				 Assemble Contig/ scaffold

What are Copy Number Variations?



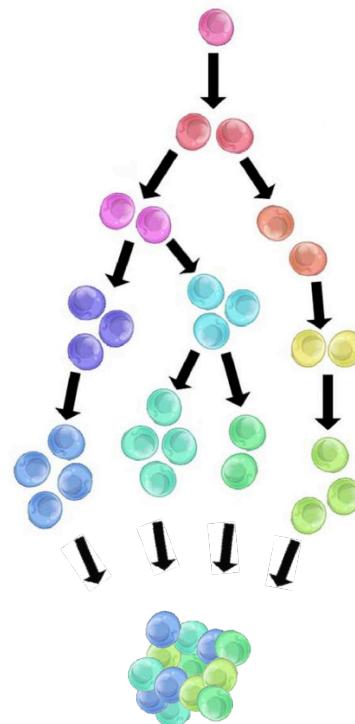
Single Cell CNV analysis

Germ Cells



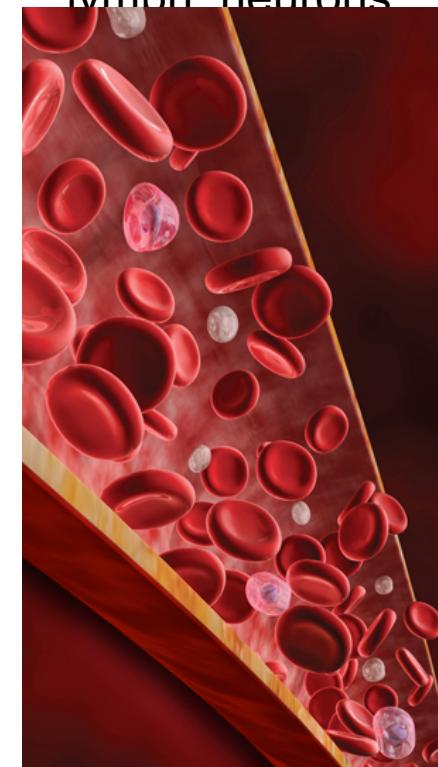
Recombination
& crossover
events

Heterogeneous
Tumors



Clonal
expansion

Heterogeneous
Tissues: blood,
lymph, neurons



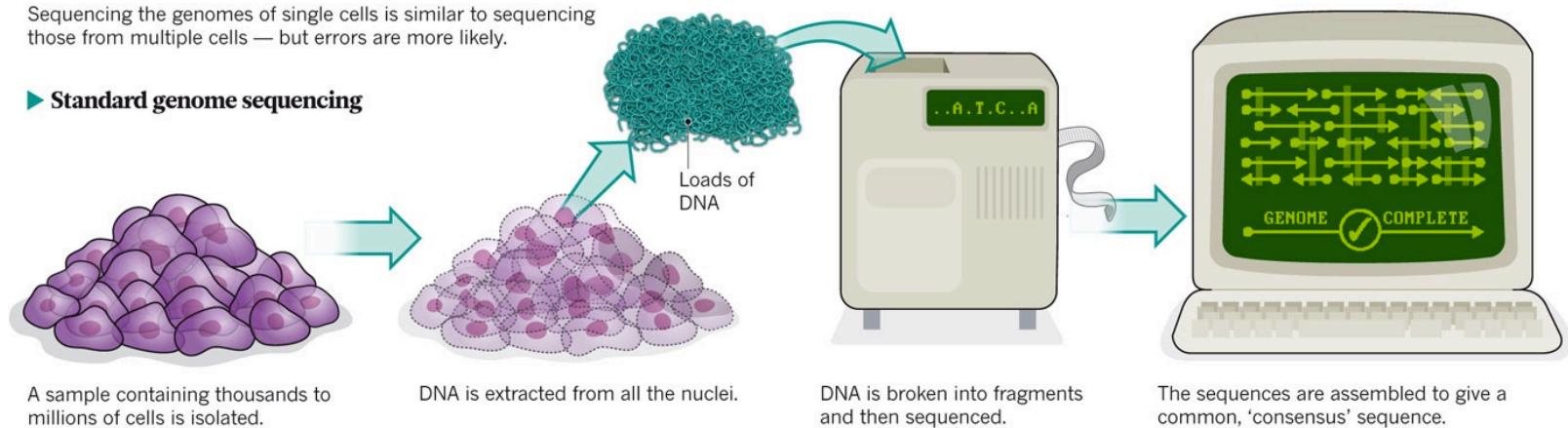
Isolating unique
cell types

Whole Genome Amplification

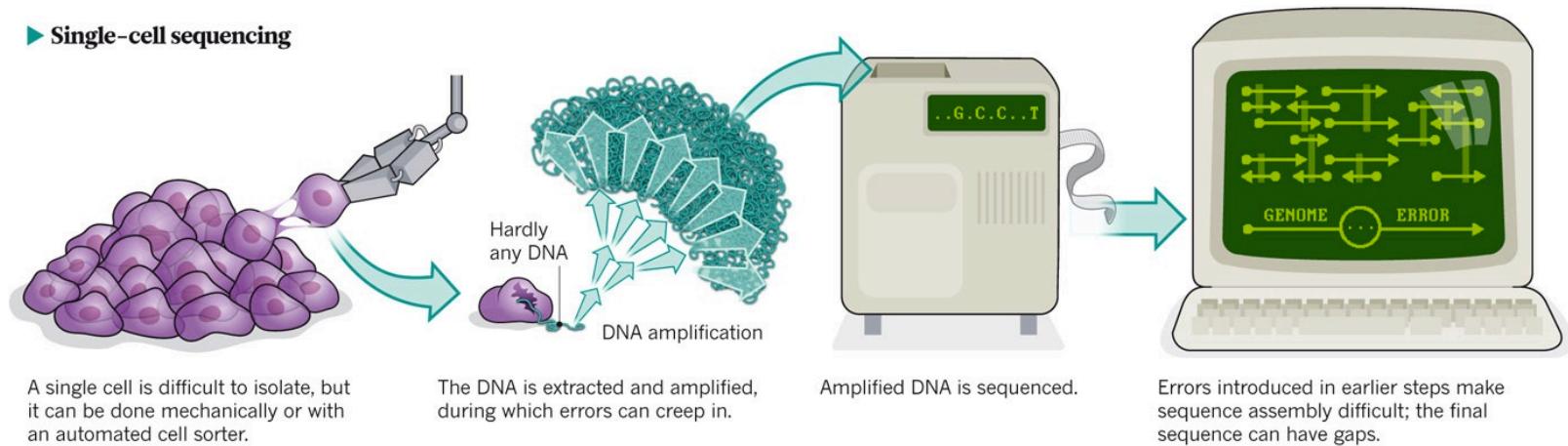
ONE GENOME FROM MANY

Sequencing the genomes of single cells is similar to sequencing those from multiple cells — but errors are more likely.

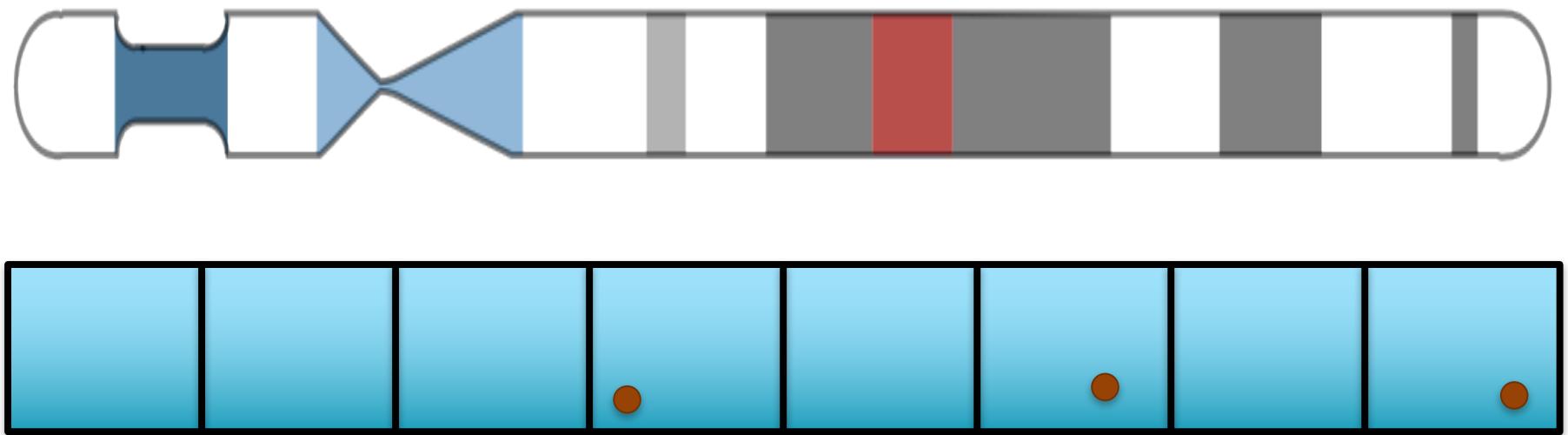
► Standard genome sequencing



► Single-cell sequencing

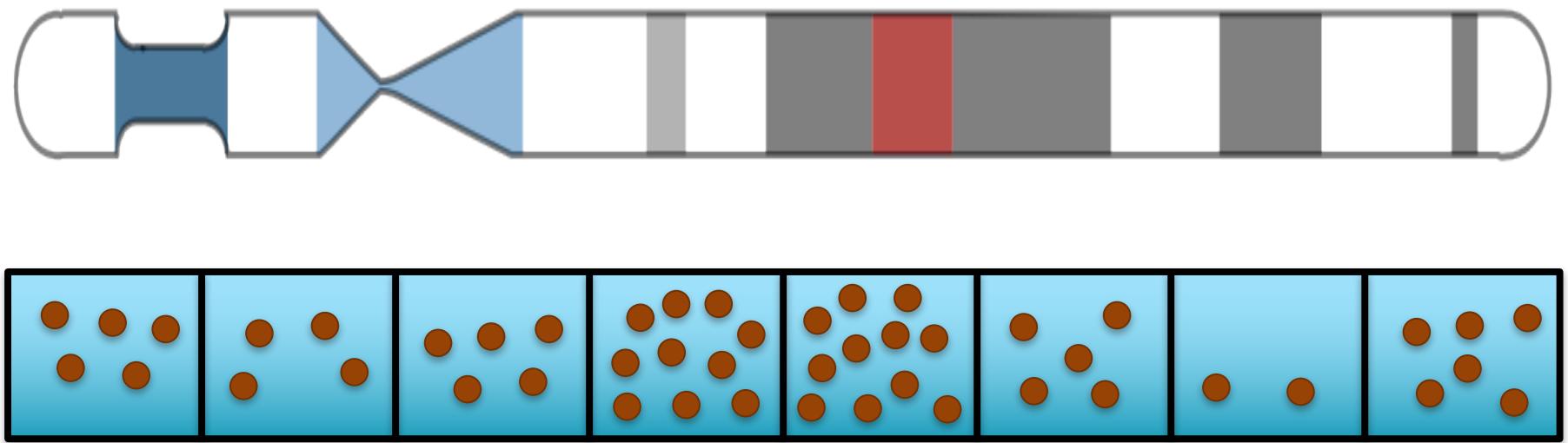


I) Binning



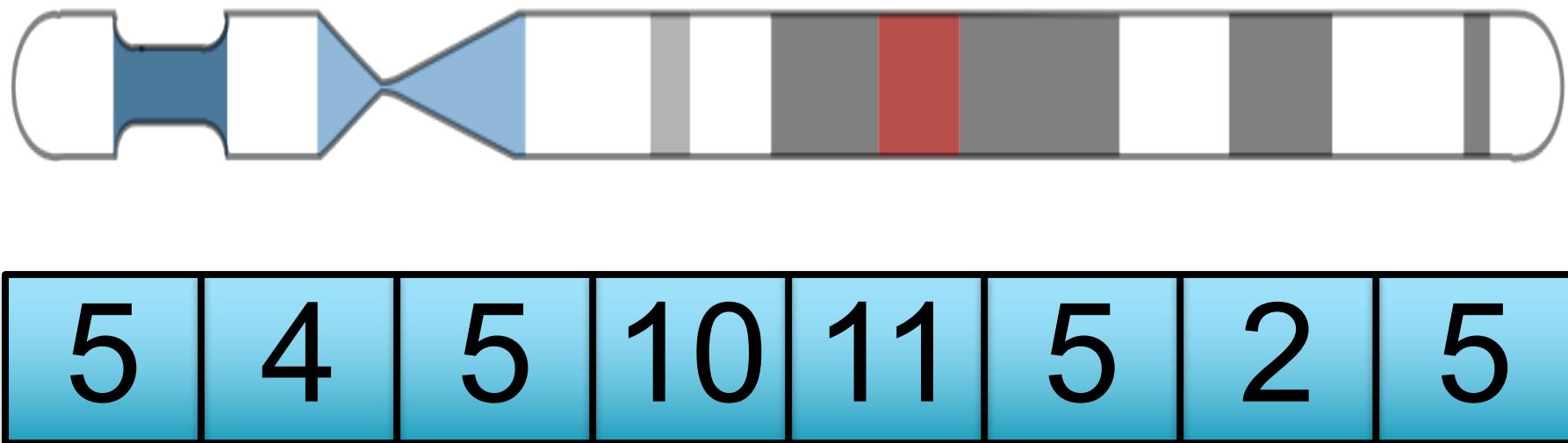
- Single Cell CNV analysis
 - Coverage is too sparse to identify point mutations
 - Divide the genome into “bins” with ~50 – 100 reads / bin
 - Map the reads and count reads per bin

I) Binning



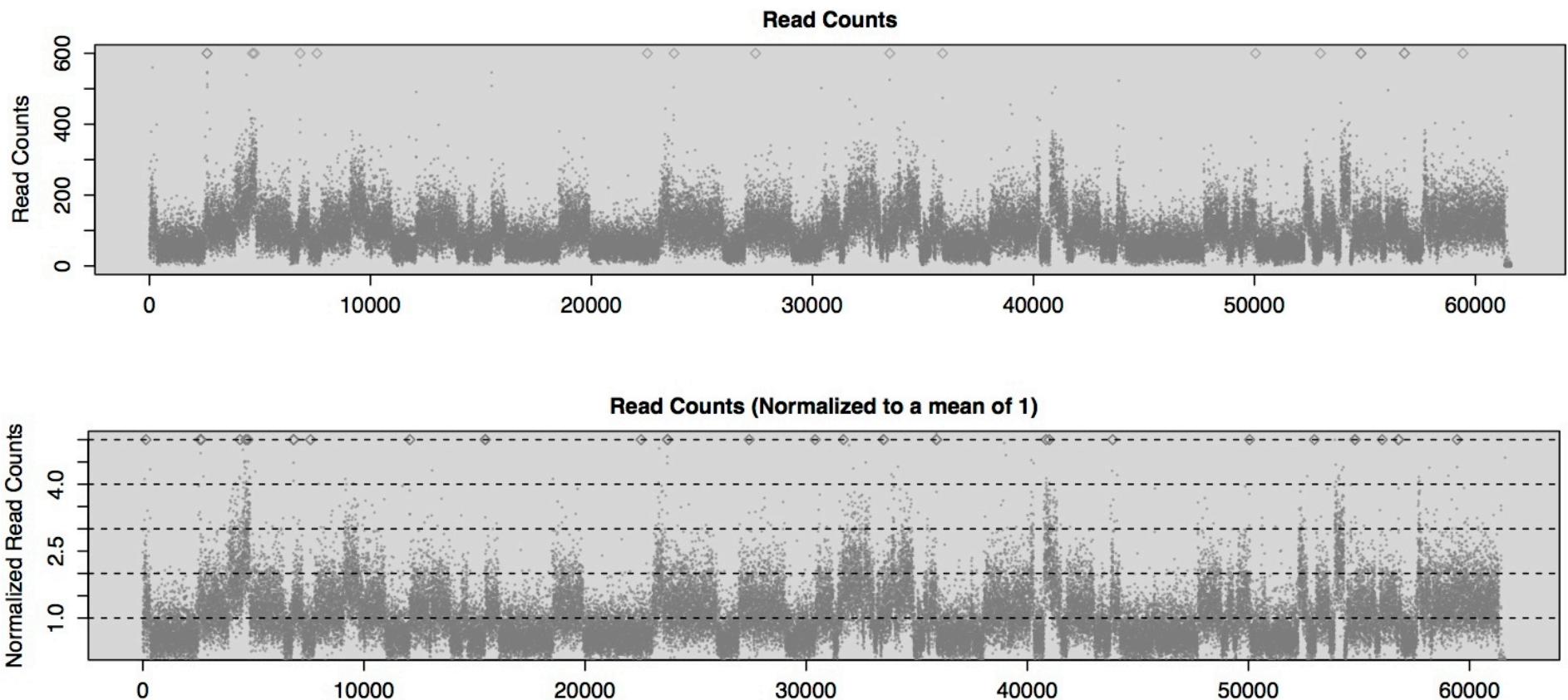
- Single Cell CNV analysis
 - Coverage is too sparse to identify point mutations
 - Divide the genome into “bins” with ~50 – 100 reads / bin
 - Map the reads and count reads per bin

I) Binning



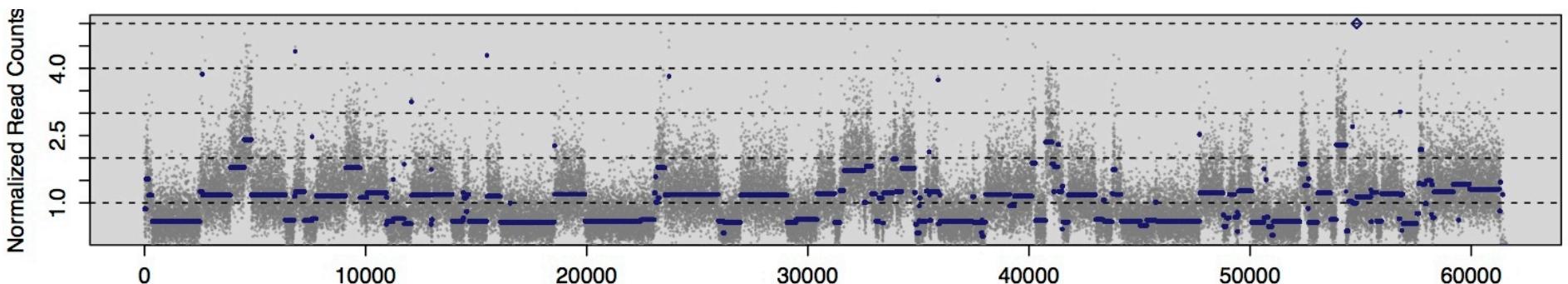
- Single Cell CNV analysis
 - Coverage is too sparse to identify point mutations
 - Divide the genome into “bins” with ~50 – 100 reads / bin
 - Map the reads and count reads per bin

2) Normalization

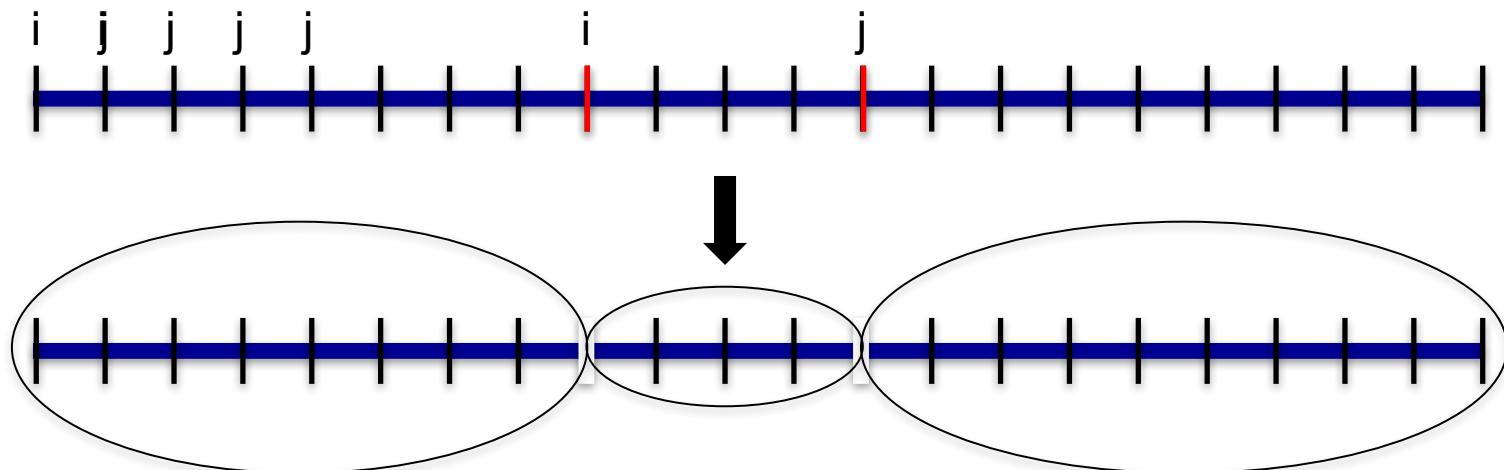


With real data, also correct for mappability, GC content, amplification biases

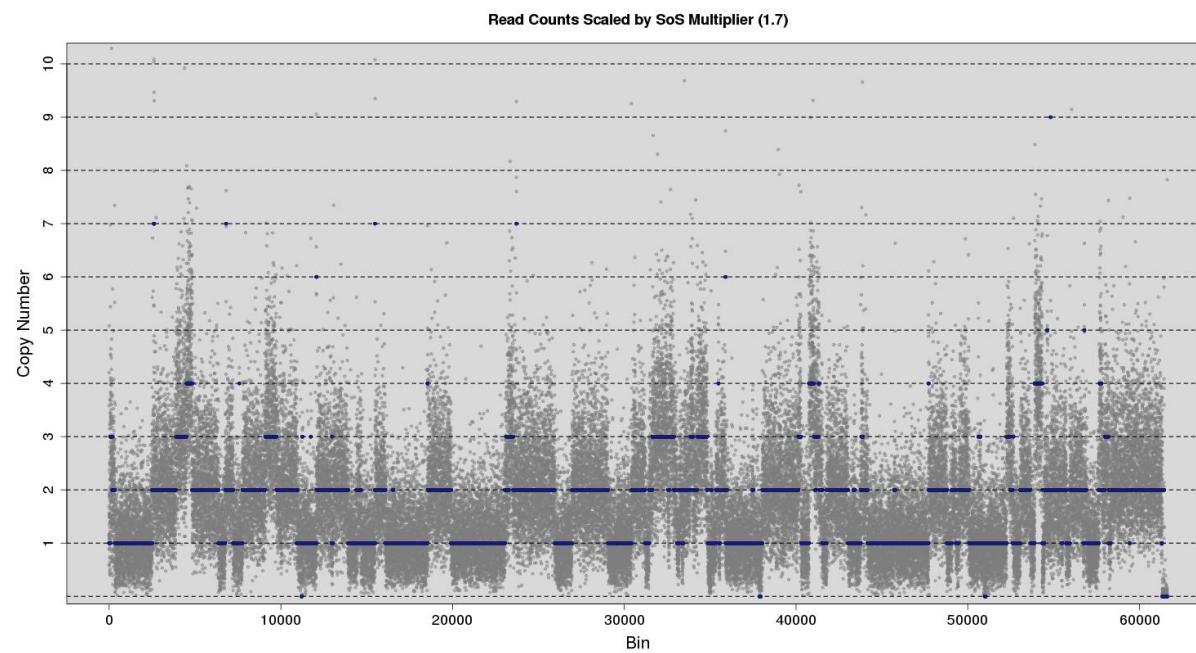
3) Segmentation



Circular Binary Segmentation (CBS)

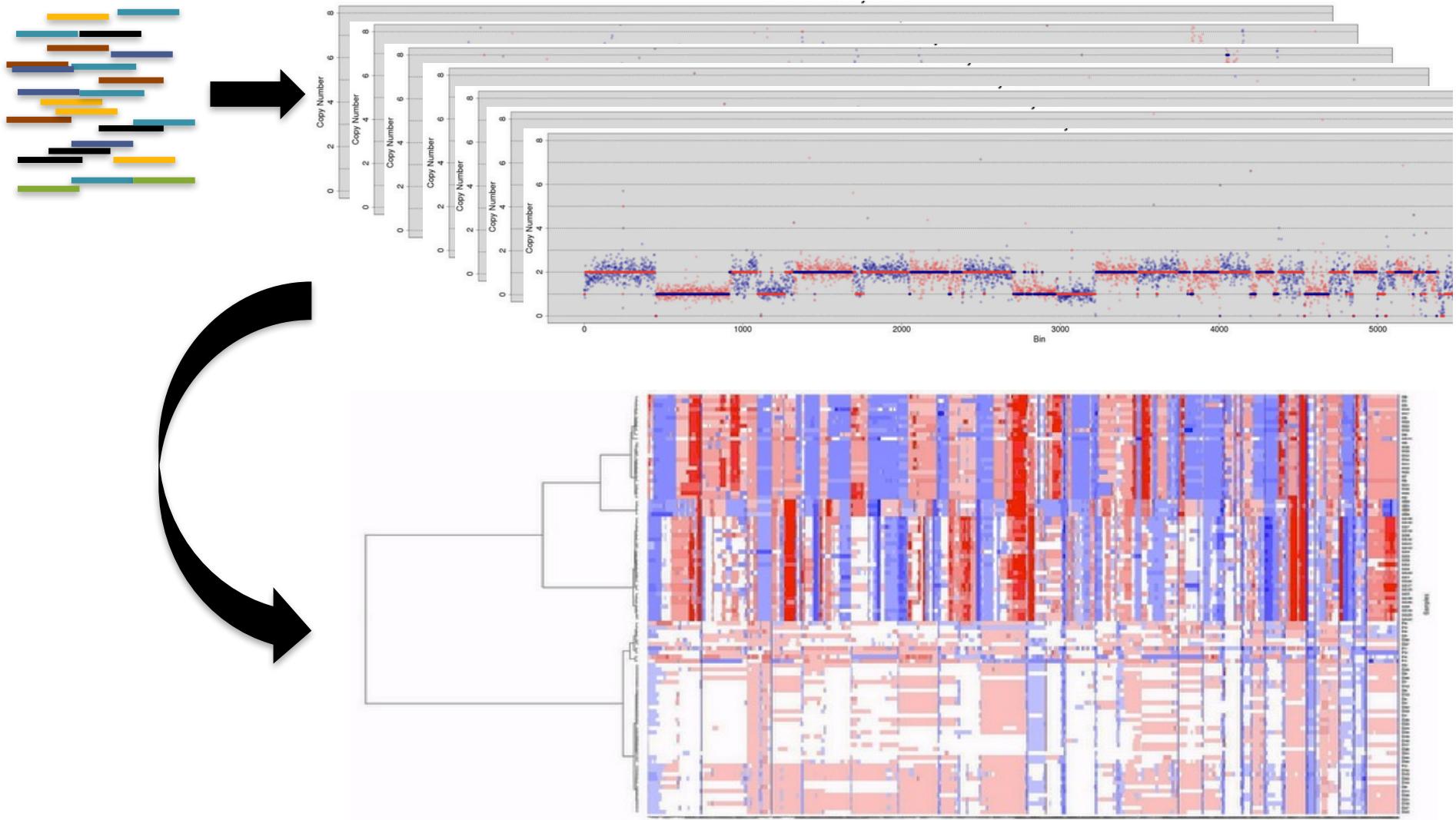


4) Estimating Copy Number



$$CN = \operatorname{argmin} \left\{ \sum_{i,j} (\hat{Y}_{i,j} - Y_{i,j}) \right\}$$

5) Cells to Populations



Interactive analysis and assessment of single-cell copy-number variations.
Garvin et al (2015) Nature Methods doi:10.1038/nmeth.3578

Structural Variation Sequence Signatures

SV classes	Read pair	Read depth	Split read	Assembly
Deletion				
Novel sequence insertion		Not applicable		
Mobile-element insertion		Not applicable		
Inversion		Not applicable		
Interspersed duplication				
Tandem duplication				



A probabilistic framework for SV discovery

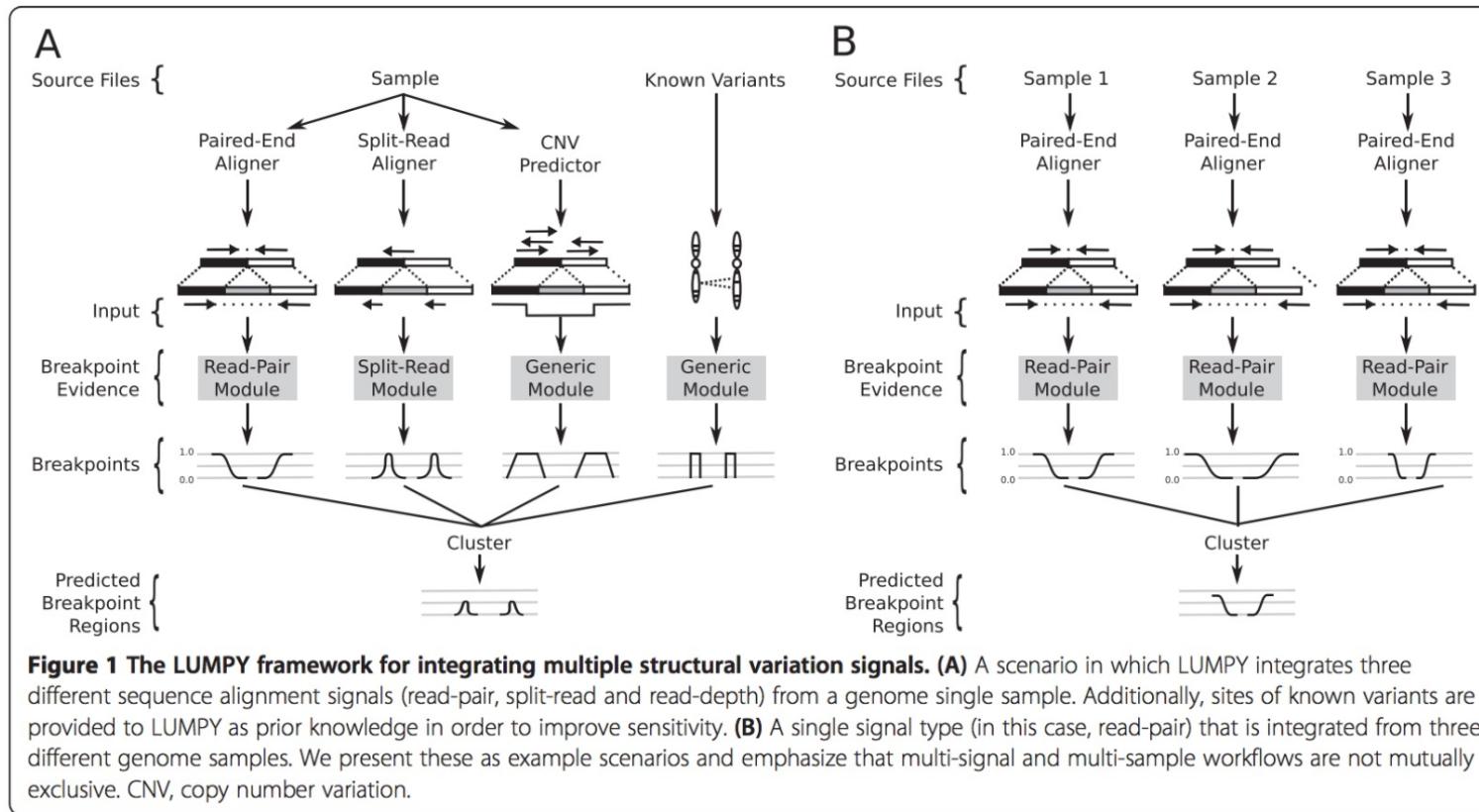


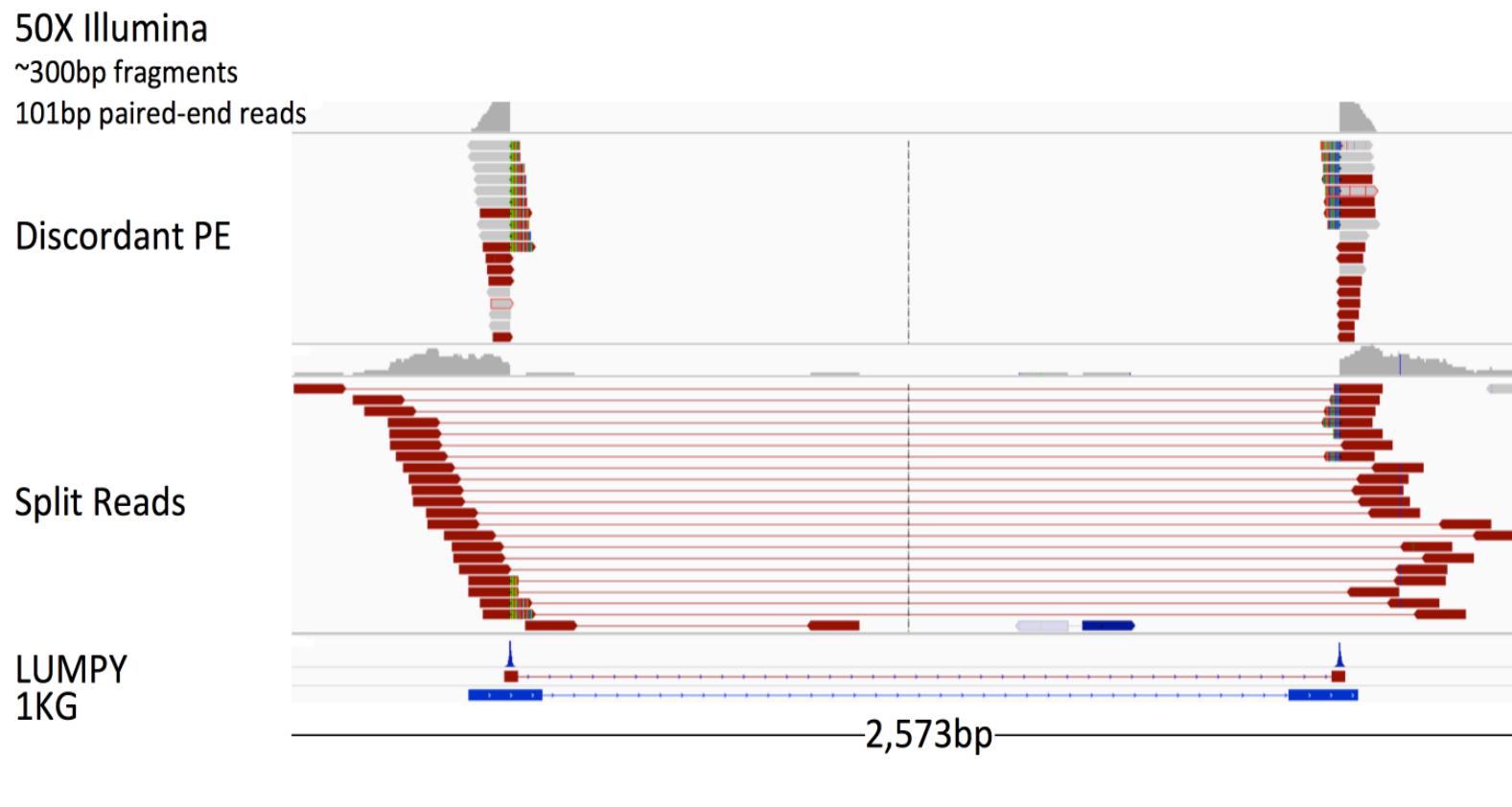
Figure 1 The LUMPY framework for integrating multiple structural variation signals. **(A)** A scenario in which LUMPY integrates three different sequence alignment signals (read-pair, split-read and read-depth) from a genome single sample. Additionally, sites of known variants are provided to LUMPY as prior knowledge in order to improve sensitivity. **(B)** A single signal type (in this case, read-pair) that is integrated from three different genome samples. We present these as example scenarios and emphasize that multi-signal and multi-sample workflows are not mutually exclusive. CNV, copy number variation.

Lumpy integrates paired-end mapping, split-read mapping, and depth of coverage for better SV discovery accuracy

Ryan Layer

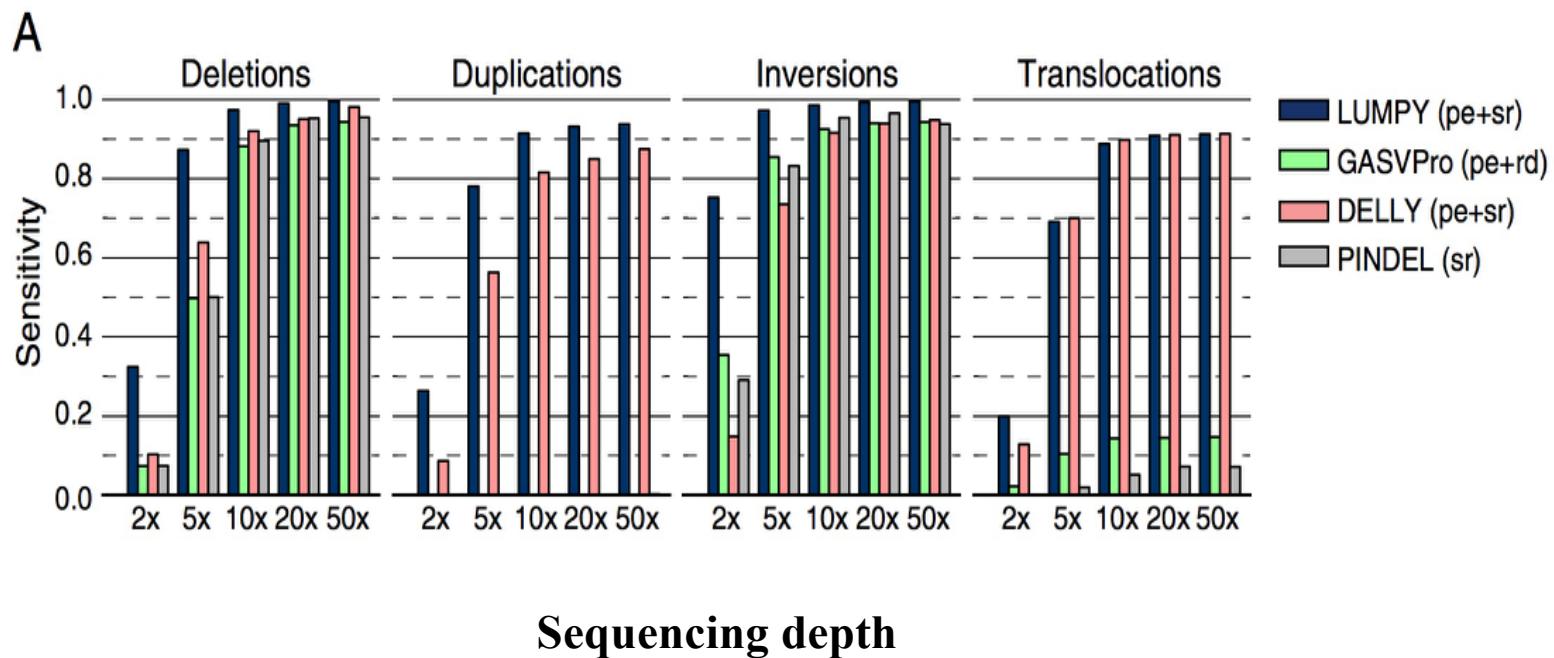


A probabilistic framework for SV discovery





A probabilistic framework for SV discovery



The dirty secrets of SV discovery

Secret #1: Often many false positives

- Short reads + heuristic alignment + rep. genome = **systematic alignment artifacts (false calls)**
- Chimeras and duplicate molecules
- Ref. genome errors (e.g., gaps, mis-assemblies)
- **ALL SV mapping studies use strict filters for above**

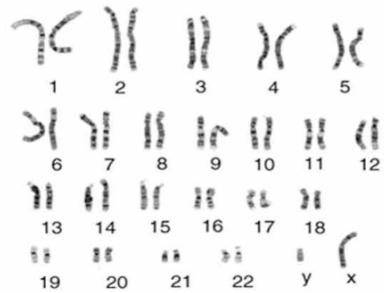
Secret #2: The false negative rate is also typically high

- Most current datasets have low to moderate *physical* coverage due to small insert size (~10-20X)
- Breakpoints are enriched in repetitive genomic regions that pose problems for sensitive read alignment
- FILTERING!
- The false negative rate is usually hard to measure, but is thought to be extremely high for most paired-end mapping studies (>30%)
- When searching for spontaneous mutations in a family or a tumor/normal comparison, a false negative call in one sample can be a false positive somatic or de novo call in another.

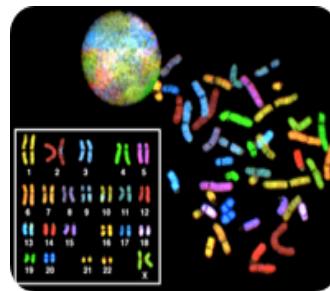
PacBio Single Molecule Real Time Sequencing (SMRT-sequencing)



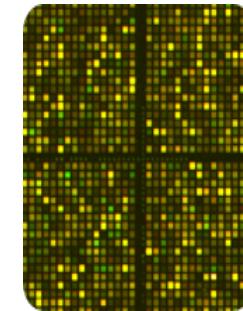
Our understanding of structural variation is driven by technology



1940s - 1980s
Cytogenetics / Karyotyping



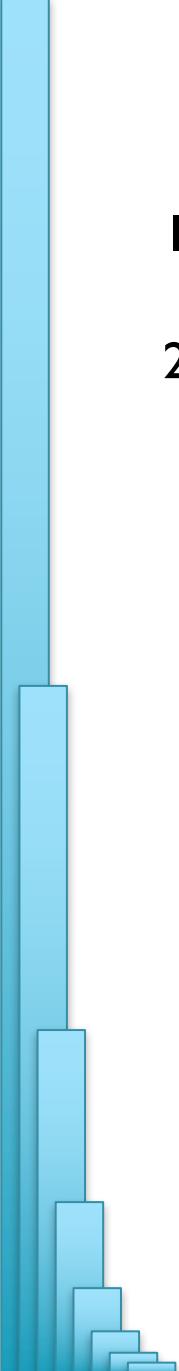
1990s
CGH / FISH /
SKY / COBRA



2000s
Genomic microarrays
BAC-aCGH / oligo-aCGH



Today
High throughput
DNA sequencing



Next Steps

- I. Finish Assignment I
2. Check out the course webpage



Welcome to Applied Comparative Genomics

<https://github.com/schatzlab/appliedgenomics>

Questions?