# In pursuit of perfect genome sequencing

Michael Schatz

December 7, 2017
UMD Institute for Genome Sciences

# In pursuit of perfect genome sequencing

1.  Why "Perfect"?

2.  What is "Perfect"?

3.  How will we achieve it?

4.  When will we achieve it?

# In pursuit of perfect genome sequencing

1. **Why "Perfect"?**

2. What is "Perfect"?

3. How will we achieve it?

4. When will we achieve it?

# The most wondrous map…



*"Without a doubt, this is the most important, most wondrous map ever produced by humankind."*

*Bill Clinton*
*June 26, 2000*

# The most wondrous map…



*"Without a doubt, this is the most important, most wondrous map ever produced by humankind."*

*Bill Clinton*
*June 26, 2000*

# Who is the reference human?

*Pieter de Jong, RPCI*

---

**Appendix: Identifying the ancestry of segments of the human genome reference sequence**
To compare Neandertal to present-day human haplotypes for the purpose of population genetic analysis, we needed to have long haploid sequences from present-day humans that were of known ancestry. To identify such segments, we took advantage of the fact that the human reference sequence is haploid over scales of tens of kilobases, because it is comprised of a tiling-path of Bacterial Artificial Chromosomes (BACs) or other clone types that are of typical size 50-150 kb (S92). We do not know of any other substantial source of high quality human haploid sequences of the requisite size.

*Determining the ancestries of the libraries in the human genome reference sequence using HAPMIX*
It is crucial to know the 'ancestry' of a clone to use it in a meaningful population genetic analysis. In what follows, we define 'ancestry' as the geographic region in which a clone's ancestor lived 1,000 years ago, inferred based on its genetic proximity to other individuals from that region today. This definition allows us to classify clones from Chinese Americans as "East Asian," from European Americans as "European", and from African Americans as either "West African" or "European".

To identify the ancestries of the libraries comprising most of the human genome reference sequence, we used a list of 26,558 clones tiling the great majority of the genome, most of which we were able to assign to a library of origin. Restricting to the autosomes, we identified 21,156 clones that seemed to fall into 9 libraries based on the naming scheme: CTA (n=199), CTB (n=356), CTC (n=452), CTD (n=1,426), RPCI-1 (n=740), RPCI-3 (n=456), RPCI-4 (n=716), RPCI-5 (n=802) and RPCI-11 (n=16,009). (In a subsequent re-examination, we identified additional clones that we likely could have classified into libraries, including 953 from RPCI-11, 632 from RPCI-1, and 490 from another library RPCI-13.) The median span of the 21,156 clones we analyzed was 112 kb, and 80% are >50kb in size. About 2/3 came from a single library, RPCI-11.

1. RPCI-11 is an African American: RPCI-11, the individual who contributed most of the human genome reference sequence, is consistent with having African American ancestry, with 42% of the clones of confident West African ancestry and 42% of the clones of confident European ancestry, and the ancestry of the remaining clones less confidently inferred. The finding of likely African American ancestry for RPCI-11 was previously reported in a study of the ancestry of RPCI-11 clones spanning the Duffy blood group locus (S93), and here we confirm this finding, and also expand the inference to the whole genome.

2. CTD is an East Asian: The majority of clones from CTD, the second largest library in its contribution to the human genome sequence, is likely an East Asian. In a HAPMIX analysis with CEU (European) – CHB+JPT (East Asian) as the proposed ancestral populations, the majority of clones are of confident East Asian origin, and there is no secondary mode of confident European ancestry, as might be expected from a Latino or South Asian individual.

3. The remaining 7 libraries are European: The remaining libraries (CTA, CTB, CTC, RPCI-1, RPCI-3, RPCI-4 and RPCI-5) are inferred to be of European ancestry, since they all have consistent distributions of inferred clone ancestries, with the majority of clones of confident European ancestry in both our HAPMIX analyses and no secondary modes.

***A Draft Sequence of the Neandertal Genome***
Green et al (2010) Science. DOI: 10.1126/science.1188021
Supplemental Note 16 (pg 145-146)

# Who is the reference human?
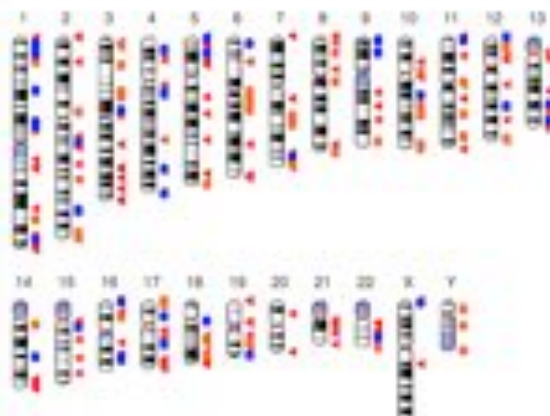
# Human Genome Overview

Information about the continuing improvement of the human genome



◀ Region containing alternate loci
● Region containing fix patches
● Region containing novel patches

Idiogram of the latest human assembly; GRCh38.p11

The GRC is working hard to provide the best possible reference assembly for human. We do this by both generating multiple representations (alternate loci) for regions that are too complex to be represented by a single path. Additionally, we are releasing regional fixes known as patches. This allows users who are interested in a specific locus to get an improved representation without affecting users who need chromosome coordinate stability.

## Download data:

- GRCh38.p11 (latest minor release) FTP
- GRCh38 (latest major release) FTP
- Genomic regions under review FTP
- Current Tiling Path Files (TPFs)

Transitioning to GRCh38? Try the NCBI Remapping Service, which uses the same assembly-assembly alignments used by the GRC.

### Next assembly update
The next assembly update (GRCh38.p12) will be a minor (patch) release in winter 2017.

---

| GRCh38.p11 | GRCh37.p13 | GRCh37 |

## GRCh38.p11

**Release date:** June 14, 2017
**Release type:** minor
**Release notes:** GRCh38.p11 is the eleventh patch release for the GRCh38 reference assembly. No chromosome coordinates changed. This release includes 11 FIX patches and 10 NOVEL patch. The total number of patch scaffolds is now: 64 FIX and 59 NOVEL.
**Assembly accessions:** GenBank: GCA_000001405.26 , RefSeq: GCF_000001405.37

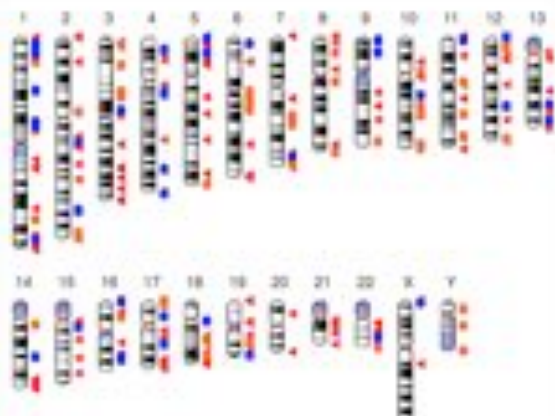Pseudoautosomal regions

| Name | Chr | Start | Stop |
|------|-----|-------|------|
| PAR#1 | X | 10,001 | 2,781,479 |
| PAR#2 | X | 155,701,383 | 156,030,895 |
| PAR#1 | Y | 10,001 | 2,781,479 |
| PAR#2 | Y | 56,887,903 | 57,217,415 |

# Genome Reference Consortium

## Human Genome Overview

Information about the continuing improvement of the human genome

The GRC is working hard to provide the best possib[...]
by both generating multiple representations (alterna[...]
represented by a single path. Additionally, we are re[...]
allows users who are interested in a specific locus t[...]
affecting users who need chromosome coordinate s[...]

### Download data:

- GRCh38.p11 (latest minor release) FTP
- GRCh38 (latest major release) FTP
- Genomic regions under review FTP
- Current Tiling Path Files (TPFs)

Transitioning to GRCh38? Try the NCBI Remapp[...]
assembly alignments used by the GRC.

**Next assembly update**
The next assembly update (GRCh38.p12) will be [...]

◄ Region containing alternate loci
● Region containing fix patches
● Region containing novel patches

Idiogram of the latest human assembly, GRCh38.p11

| GRCh38.p11 | GRCh37.p13 | GRCh37 |

## GRCh38.p11

**Release date:** June 14, 2017
**Release type:** minor
**Release notes:** GRCh38.p11 is the eleventh patch release for the GRCh38 reference assembly. No chromosome coordin[...]
of patch scaffolds is now: 64 FIX and 59 NOVEL.
**Assembly accessions:** GenBank: GCA_000001405.26 , RefSeq: GCF_000001405.37

Pseudoautosomal regions

| Name | Chr | Start | Stop |
|------|-----|-------|------|
| PAR#1 | X | 10,001 | 2,781,479 |
| PAR#2 | X | 155,701,383 | 156,030,895 |
| PAR#1 | Y | 10,001 | 2,781,479 |
| PAR#2 | Y | 56,887,903 | 57,217,415 |

# Importance of Personal Genomes

**Current standard is to align your data to the "reference" human genome.**

**But the "reference" isn't really the genome for *any* human and potentially biases the results in many ways:**

- ***Genome:*** biased read mapping, causing false positive and false negative mutations

- ***Transcriptome***: mutations of splice sites, stop codons or branch point change gene models, CNVs modulate expression levels, gene fusions create new transcripts

- ***Epigenome***: *cis* versus *trans* effects, *allele-specific expression, allele-specific binding*

> ***Same issues apply to most "reference" genomes***

# In pursuit of perfect genome sequencing

1. **Why "Perfect"?**

   *Because it is important, complex, and personal*

2. What is "Perfect"?

3. How will we achieve it?

4. When will we achieve it?

# In pursuit of perfect genome sequencing

1. Why "Perfect"?

2. **What is "Perfect"?**

3. How will we achieve it?
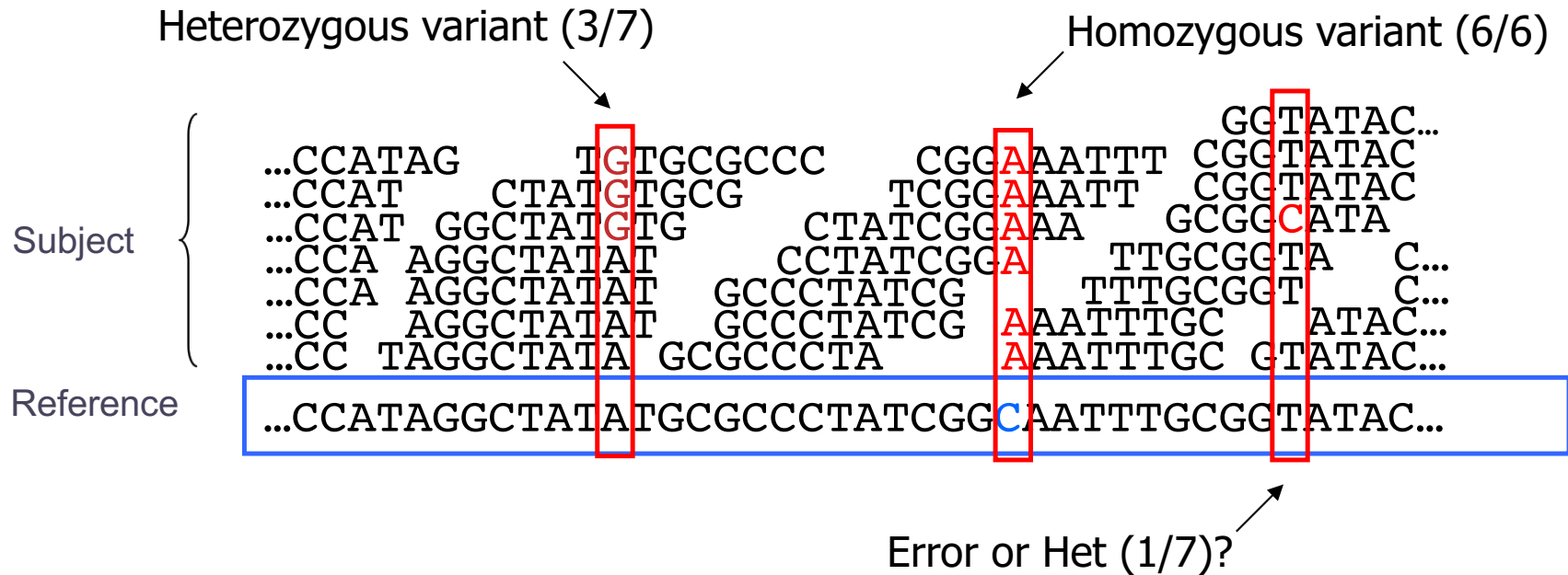
4. When will we achieve it?

# *1. Correctness*:

Is the genome faithfully represented?
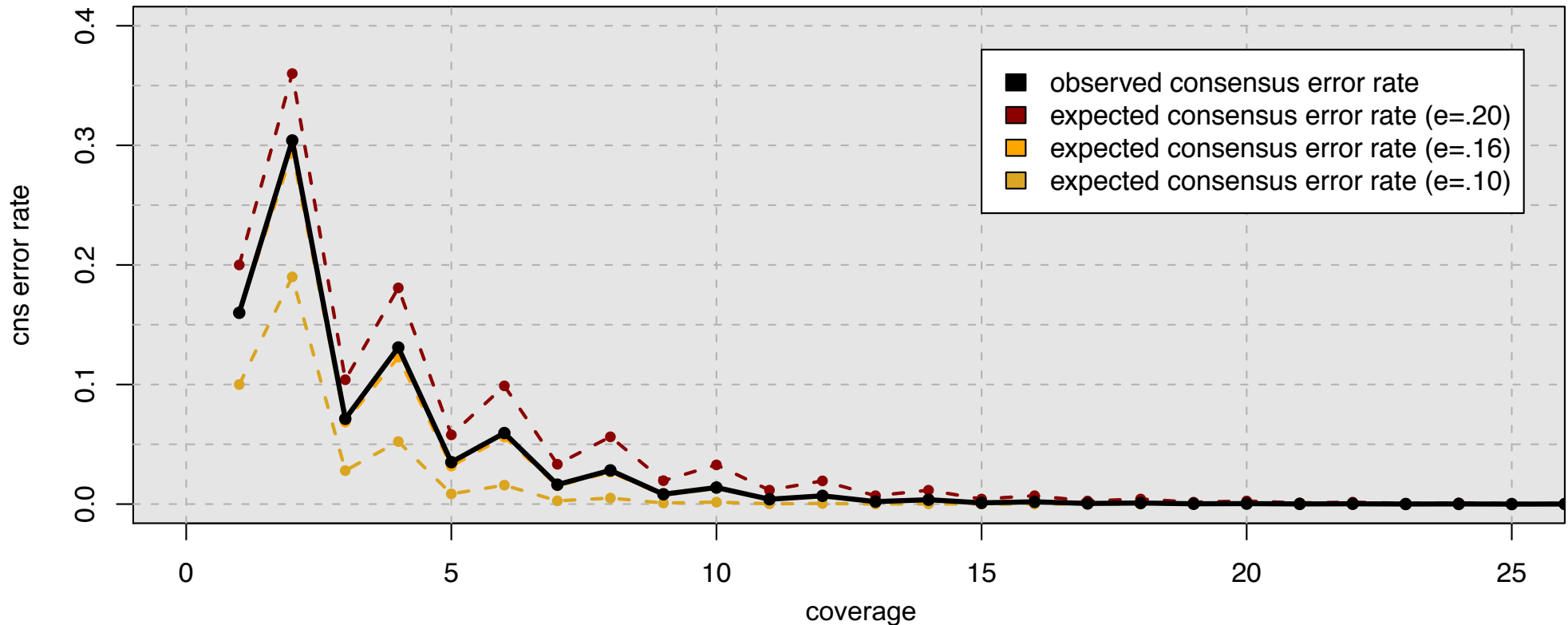
# 1. Correctness:
## Is the genome faithfully represented?



PacBio RS II

CSHL/PacBio

```
TTGTAAGCAGTTGAAAACTATGTGTGGATTTAGAATAAAGAACATGAAAG
||||||||||||||||||||||||| ||||||| |||||||||||| |||
TTGTAAGCAGTTGAAAACTATGTGT-GATTTAG-ATAAAGAACATGGAAG

ATTATAAA-CAGTTGATCCATT-AGAAGA-AAACGCAAAAGGCGGCTAGG
| |||||| ||||||||||||| |||| | |||||| |||||| ||||||
A-TATAAATCAGTTGATCCATTAAGAA-AGAAACGC-AAAGGC-GCTAGG

CAACCTTGAATGTAATCGCACTTGAAGAACAAGATTTTATTCCGCGCCCG
| |||||| |||| |||| || |||||||||||||||||||||||||||||
C-ACCTTG-ATGT-AT--CACTTGAAGAACAAGATTTTATTCCGCGCCCG

TAACGAATCAAGATTCTGAAAACACAT-ATAACAACCTCCAAAA-CACAA
| ||||||||| |||||||||||||| || || |||||||||| ||||||
T-ACGAATC-AGATTCTGAAAACA-ATGAT----ACCTCCAAAAGCACAA

-AGGAGGGGAAAGGGGGGAATATCT-ATAAAAGATTACAAATTAGA-TGA
 ||||||    ||     |||||||| || |||||||||||||| || |||
GAGGAGG---AA-----GAATATCTGAT-AAAGATTACAAATT-GAGTGA

ACT-AATTCACAATA-AATAACACTTTTA-ACAGAATTGAT-GGAA-GTT
||| ||||||||| | |||||||||||| ||| |||||||| |||| |||
ACTAAATTCACAA-ATAATAACACTTTTAGACAAAATTGATGGGAAGGTT

TCGGAGAGATCCAAAACAATGGGC-ATCGCCTTTGA-GTTAC-AATCAAA
|| ||||||||| ||||||| ||| ||| |||| |||||| ||| |||||||
TC-GAGAGATCC-AAACAAT-GGCGATCG-CTTTGACGTTACAAATCAAA
```

Sample of 100k reads aligned with BLASR requiring >100bp alignment
Average overall accuracy 83.7%: 11.5% insertions, 3.4% deletions, 1.4% mismatch

# Genotyping Theory



- If there were no sequencing errors, identifying SNPs would be trivial:
  - Any time a read disagrees with the reference, it must be a variant!

- A single read of many differing from the reference is probably just an error, but it becomes more likely to be real as we see it multiple times
  - Use binomial test to evaluate prob. of heterozygosity vs. prob of error
  - Coverage (oversampling) is our main tool to improve accuracy

# Consensus Accuracy and Coverage



**Coverage can overcome random errors**

- Dashed: error model from binomial sampling
- Solid: observed accuracy

$$CNS\,Error\ =\ \sum_{i=\lceil c/2 \rceil}^{c} \binom{c}{i} (e)^i (1-e)^{n-i}$$

# FALCON Accuracy



"*The overall base-to-base concordance rate is about 99.99*% (QV40 in Phred scale) in the F1 FALCON-Unzip assembly. The insertion and deletion (indel) concordances to the parental lines were lower (about QV40) than the SNP concordance rate (about QV50), with most residual errors concentrated in long homopolymer sequences"

**Phased Diploid Genome Assembly with Single Molecule Real-Time Sequencing**
Chin et al (2016) *Nature Methods. doi:10.1038/nmeth.4035.*

# 2. *Completeness*:

How much of the genome is present?

# 2. Completeness:
## How much of the genome is present?



*"88% of GWAS SNPs are intronic or intergenic of unknown function"*
ENCODE Consortium (2012) Nature

# Non-coding Somatic SNVs in PDAC



Input whole genome sequencing data
1) Matched tumor-normal SNV calls
2) RNA-seq expression calls

↓

**FunSeq2**
Prioritize non-coding regulatory variants

For each CRR variant | For each CRR class

Associate recurrently mutated CRRs with flanking genes | Determine mutation rates for each regulatory class

Use permutation testing to identify CRRs affecting expression | Normalize mutation rates for GC content, size, and abundance

Generate false discovery rates | Compute expression modulation scores

Pathway analysis
Patient survival analysis

*Coding alterations of PDAC are now fairly well established but non-coding mutations (NCMs) largely unexplored*

- Developed GECCO to analyze the thousands of somatic mutations observed from hundreds of tumors to find potential drivers of gene expression and pathogenesis

- NCMs are enriched in known and novel pathways
- NCMs correlate with changes in gene expression
- NCMs can demonstrably modulate gene expression
- NCMs correlate with novel clinical outcomes

*NCMs are an important mechanism for tumor genome evolution*

**Recurrent noncoding regulatory mutations in pancreatic ductal adenocarcinoma**
Feigin, M, Garvin, T et al. (2017) Nature Genetics. doi:10.1038/ng.3861

# Structural Variations



Deletion · Novel sequence insertion · Mobile-element insertion · Tandem duplication · Interspersed duplication · Inversion · Translocation
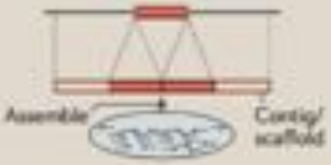
**Any mutation >50bp**

**Profound impact on genome structure and function**

# Structural Variation Sequence Signatures

# Structural Variation Sequence Signatures



| SV classes | Read pair | Read depth | Split read | Assembly |
|---|---|---|---|---|
| Deletion | | | | |

# Structural Variation Sequence Signatures



| SV classes | Read pair | Read depth | Split read | Assembly |
|---|---|---|---|---|
| Deletion | | | | Contig/scaffold Assemble |



**PacBio Sequel**



**Oxford Nanopore MinION**

***Long Read Single Molecule Sequencing***
*No Amplification Artifacts*
*Improved Mapping & De novo assemblies*
*Complete Genomes with all variant types*
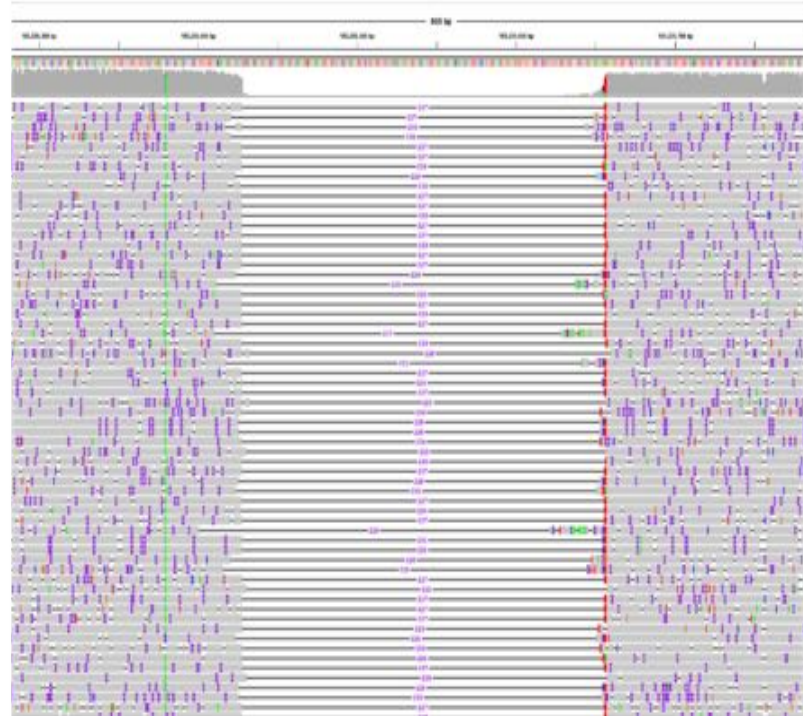
# NGMLR + Sniffles

BWA-MEM:
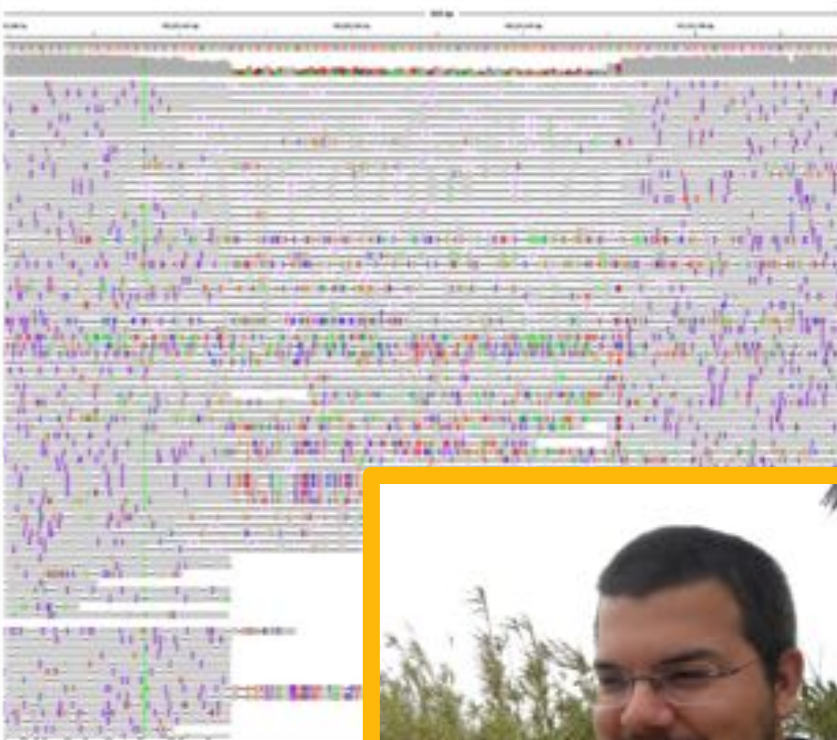


***Accurate detection of complex structural variations using single molecule sequencing***
Sedlazeck, Rescheneder et al (2017) *bioRxiv https://doi.org/10.1101/169557*

# NGMLR + Sniffles

BWA-MEM:

NGMLR:



***Accurate detection of complex structural variations using single molecule sequencing***
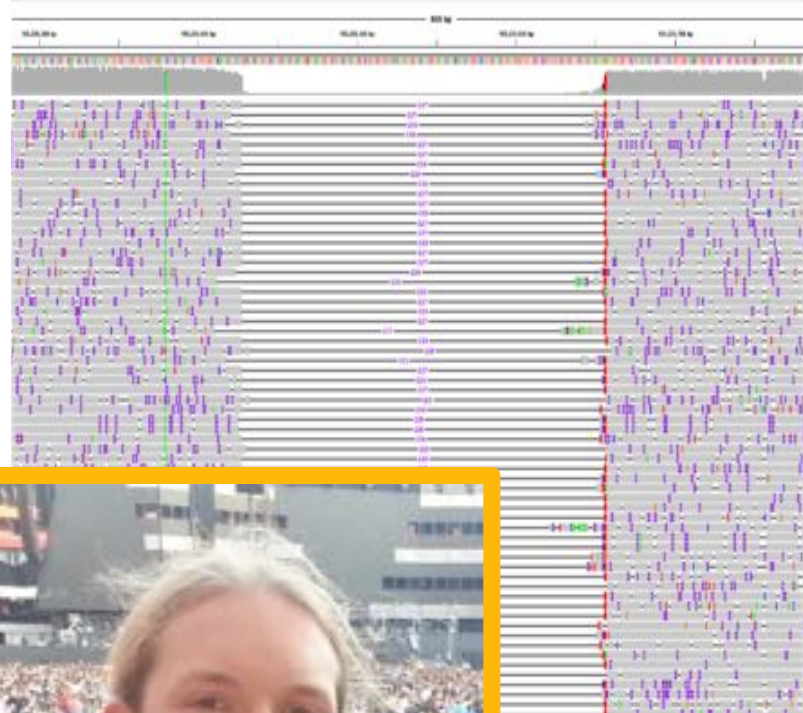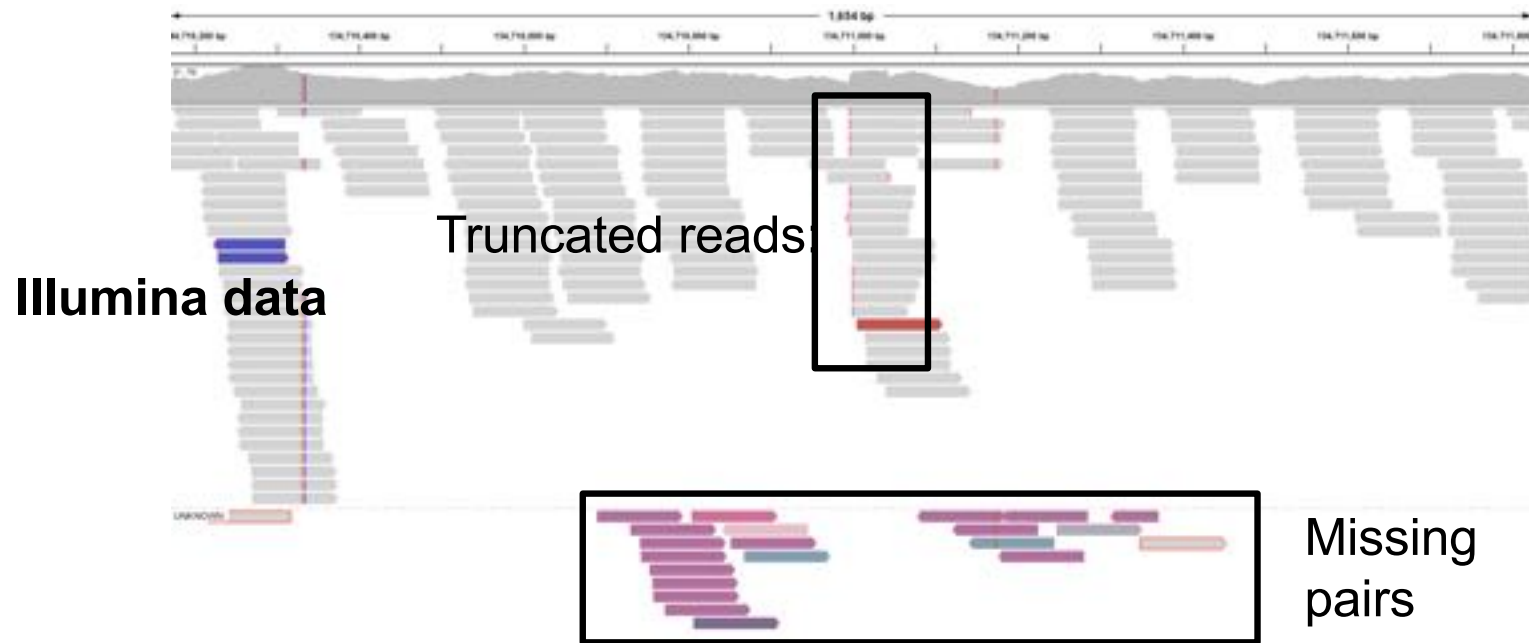Sedlazeck, Rescheneder et al (2017) *bioRxiv https://doi.org/10.1101/169557*
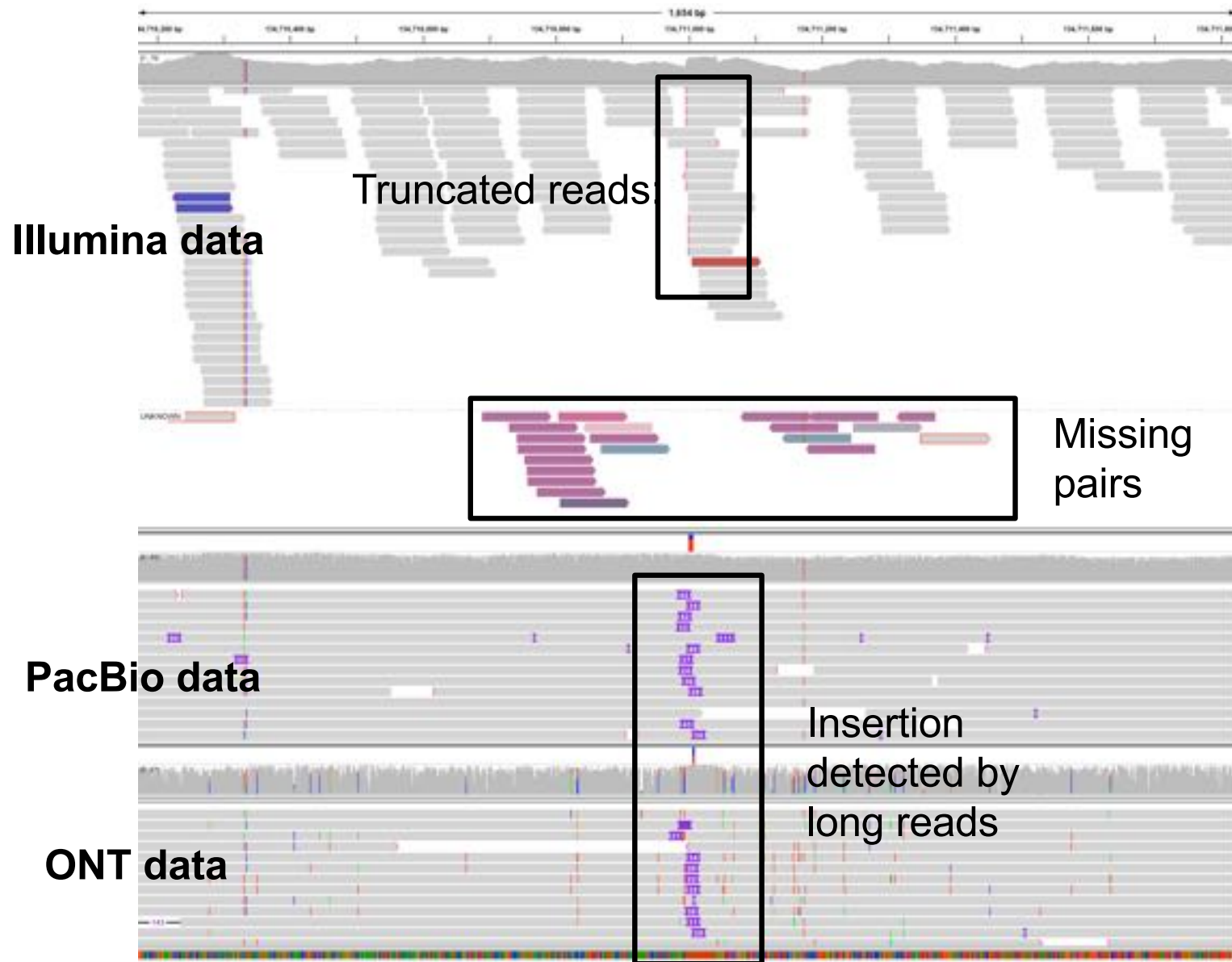
# NGMLR + Sniffles

BWA-MEM:

NGMLR:



*Accurate detection of complex structural variations using single molecule sequencing*
Sedlazeck, Rescheneder et al (2017) *bioRxiv https://doi.org/10.1101/169557*

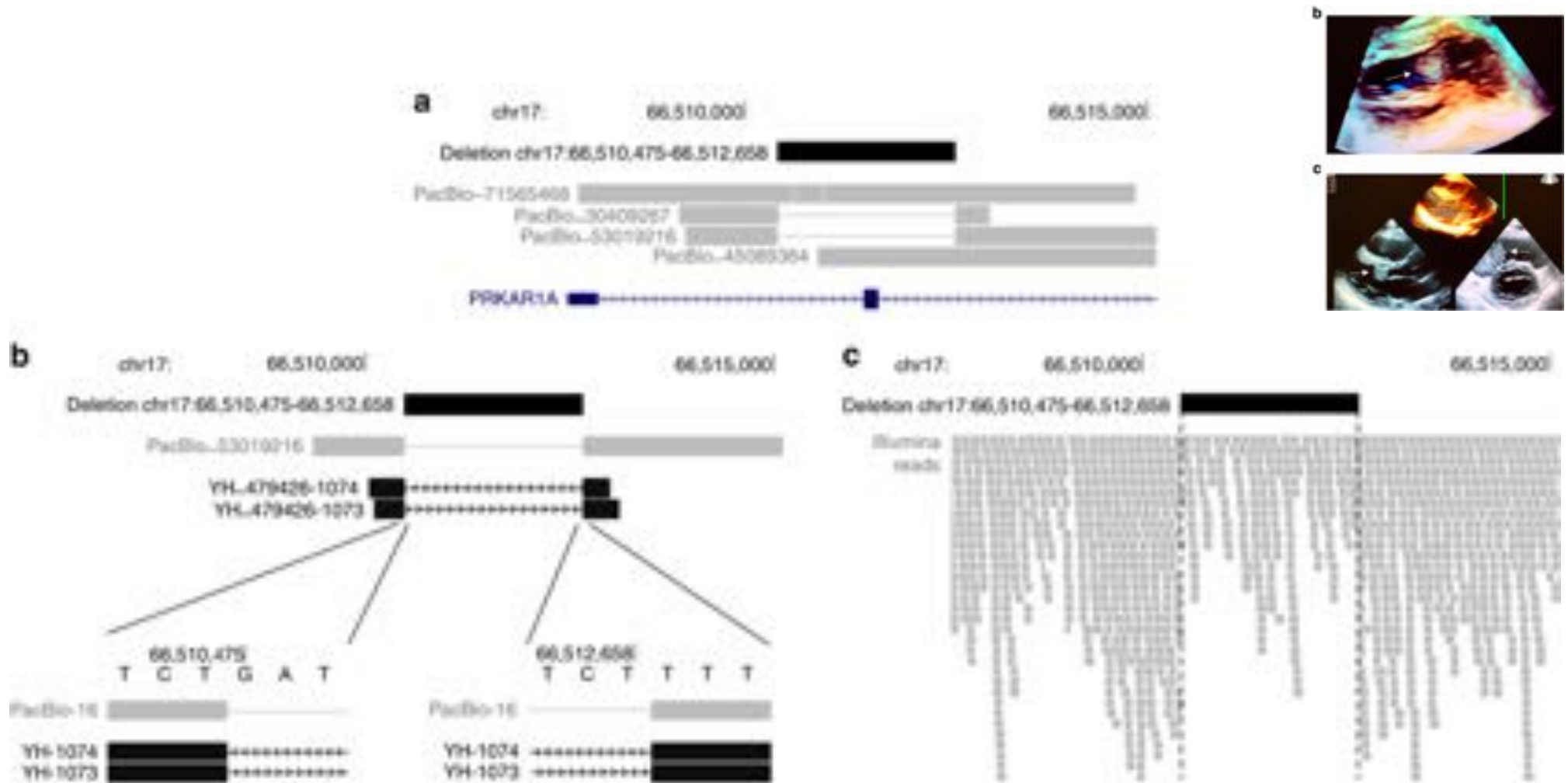# No more false positives!



Illumina data

Truncated reads

Missing pairs

*Accurate detection of complex structural variations using single molecule sequencing*
Sedlazeck, Rescheneder et al (2017) *bioRxiv https://doi.org/10.1101/169557*

# No more false positives!



Illumina data

Truncated reads

Missing pairs

PacBio data

ONT data

Insertion detected by long reads

*Accurate detection of complex structural variations using single molecule sequencing*
Sedlazeck, Rescheneder et al (2017) *bioRxiv https://doi.org/10.1101/169557*

# Structural Variations in Mendelian Disease

# Structural Variations in Breast Cancer



Figure 1 | Variants found in SK-BR-3 with PacBio long-read sequencing. (A) Circos plot showing long-range (larger than 10 kbp or interchromosomal) variants found by Sniffles from split-read alignments, with read coverage shown in the outer track. (B) Variant size histogram of deletions and insertions from size 50 bp up to 1 kbp found by log-read (Sniffles) and short-read (Survivor 2-caller consensus) variant-calling, showing similar size distributions for insertions and deletions from long reads but not for short reads where insertions are entirely missing. (C) Sniffles variant counts by type for variants above 1 kbp in size, including translocations and inverted duplications.

*Complex rearrangements and oncogene amplifications revealed by long-read DNA and RNA sequencing of a breast cancer cell line*
Nattestad, M et al (2017) bioRxiv https://doi.org/10.1101/174938

# In pursuit of perfect genome sequencing

1.  Why "Perfect"?

2.  **What is "Perfect"?**
    *100% Correct & 100% Complete*

3.  How will we achieve it?
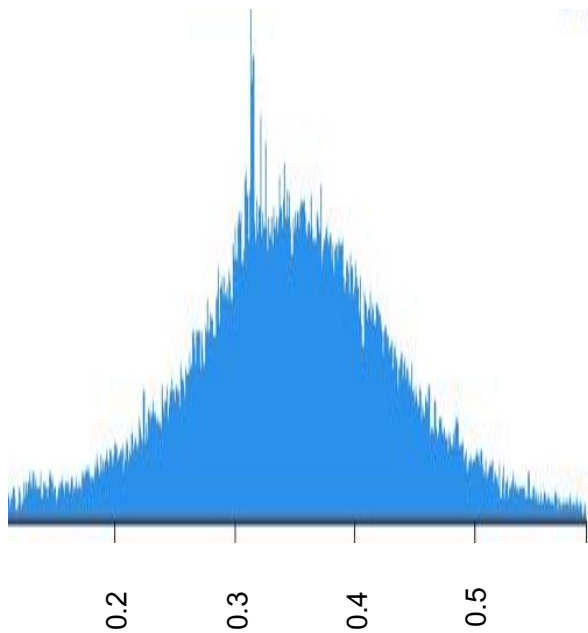
4.  When will we achieve it?

# In pursuit of perfect genome sequencing

1. Why "Perfect"?

2. What is "Perfect"?

3. **How will we achieve it?**

4. When will we achieve it?
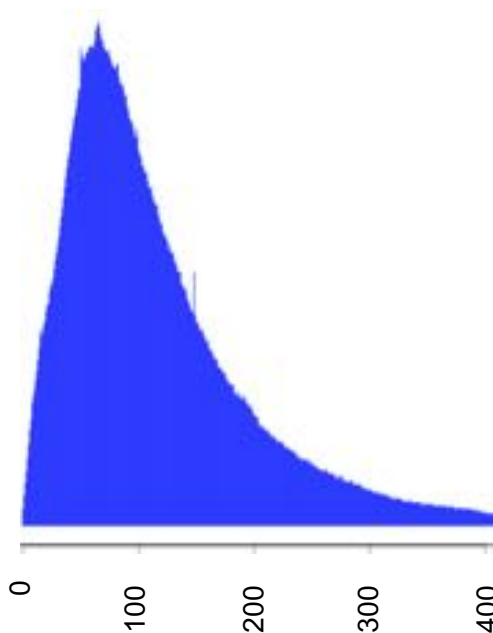
# Human Genome Sequencing Data



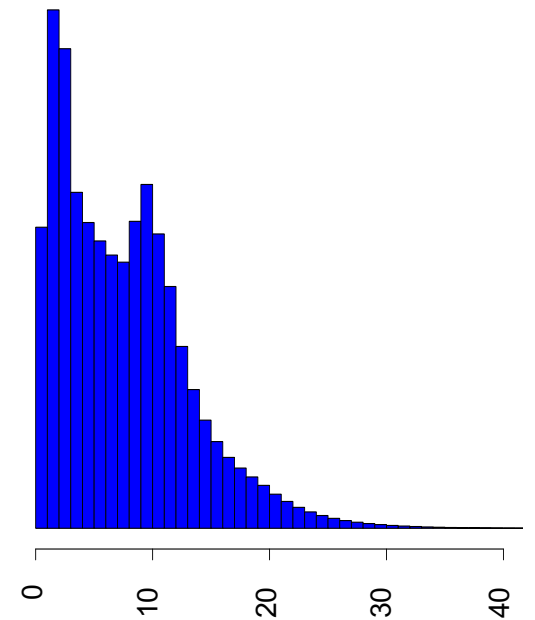Illumina — Fragment Length (kbp)
**60x Paired End**
μ=350bp

10X Genomics — Molecule Length (kbp)
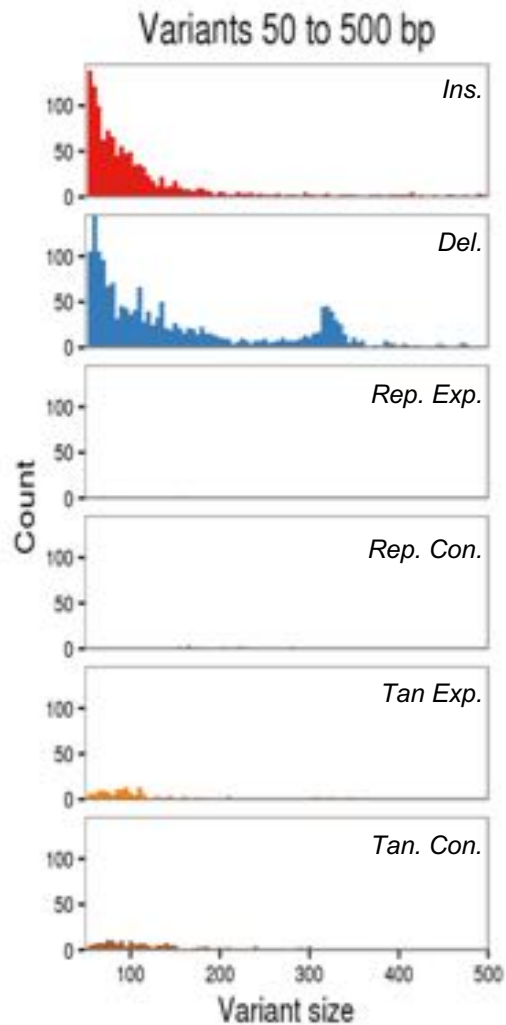**35x Linked Reads**
μ=117kbp

PacBio — Read Length (kbp)
**55x Long Reads**
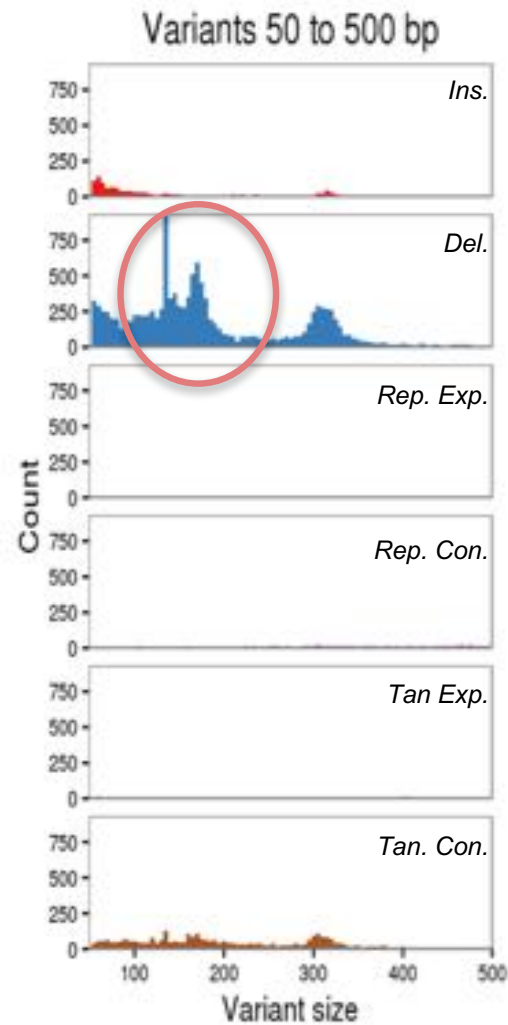μ=7.5kbp

# Missing Insertions from Short and Linked Read?

# Structural Variations Concordance



**Main Diagonal**
- Calls per tool

**Outer triplets**
- Concordance by Technology

**Inner triplets**
- Concordance by Assembly
- Concordance by Mappers

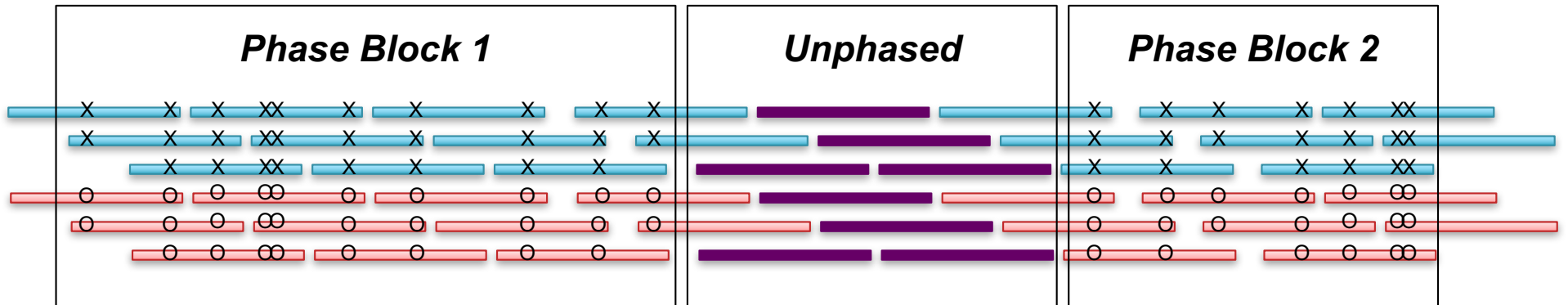**Overall:**
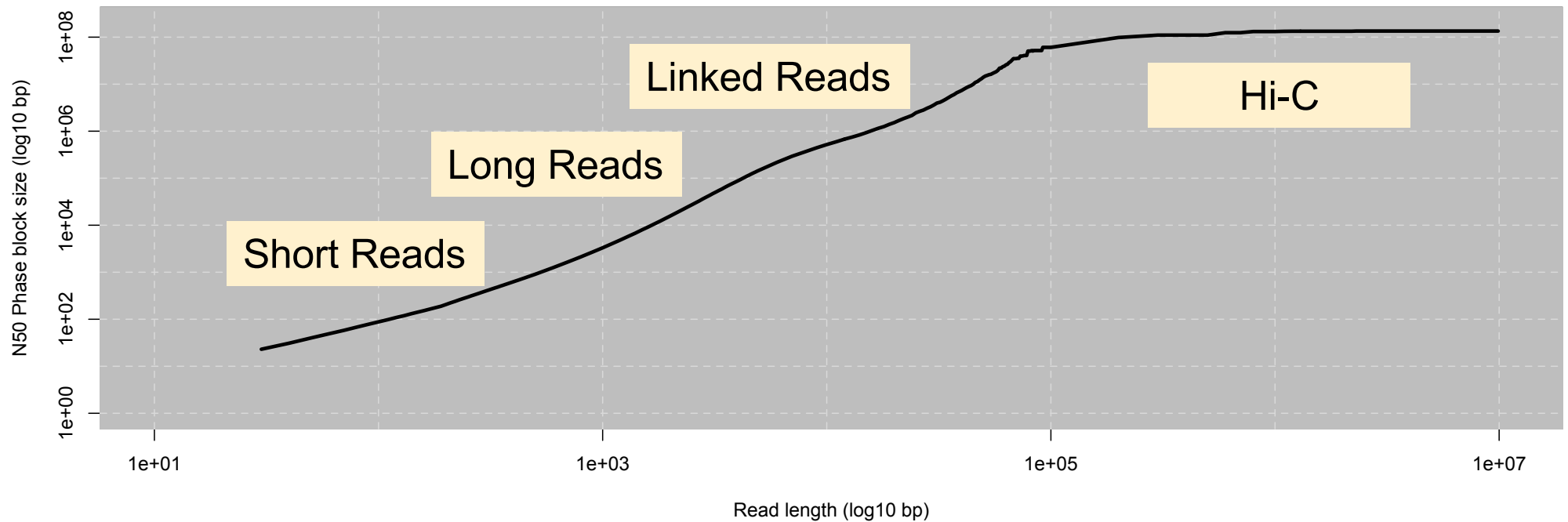- Lonnnnnnng reads give the most variants with the best concordance ☺

# Phasing Results



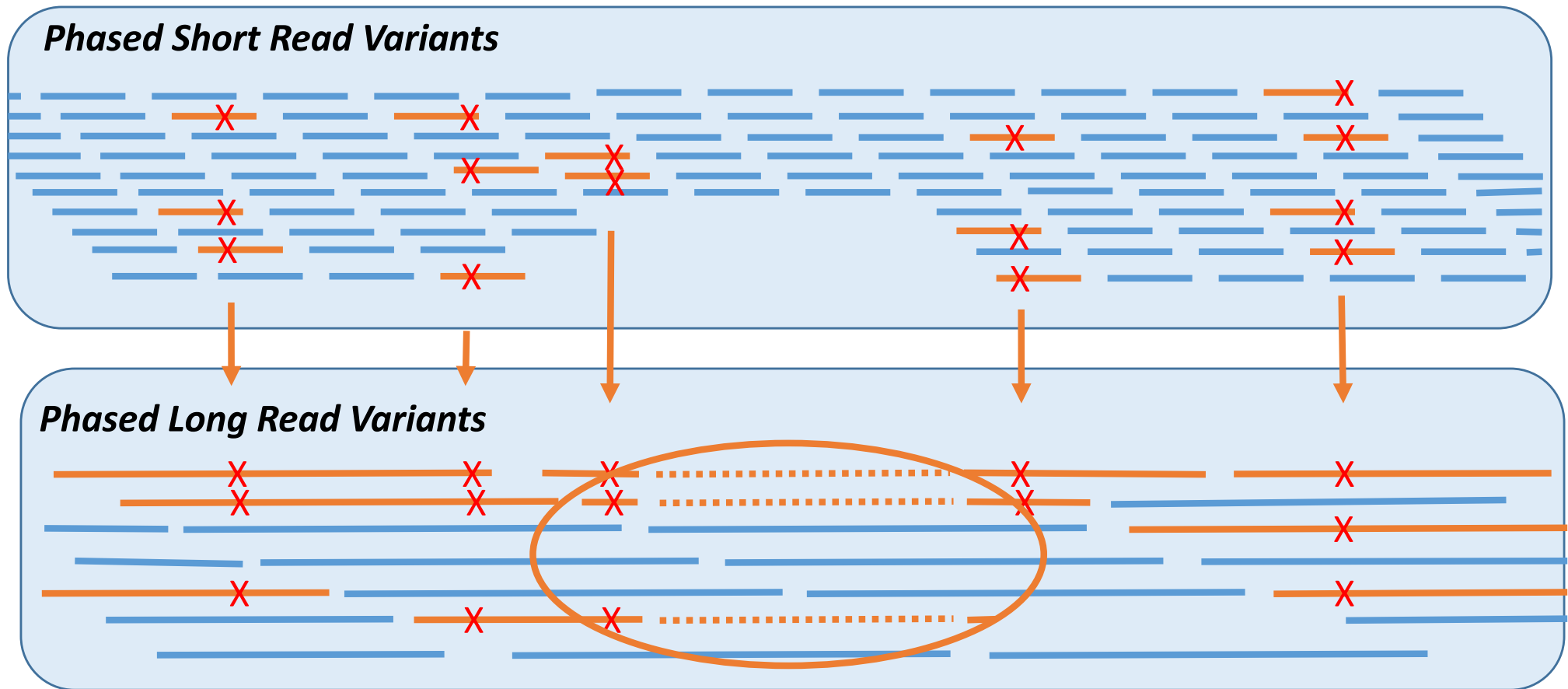**NA12878 Optimal phase block length increases with read length**



*Piercing the dark matter: Bioinformatics for third generation sequencing*
Sedlazeck *et al* (2017) *Under Review*
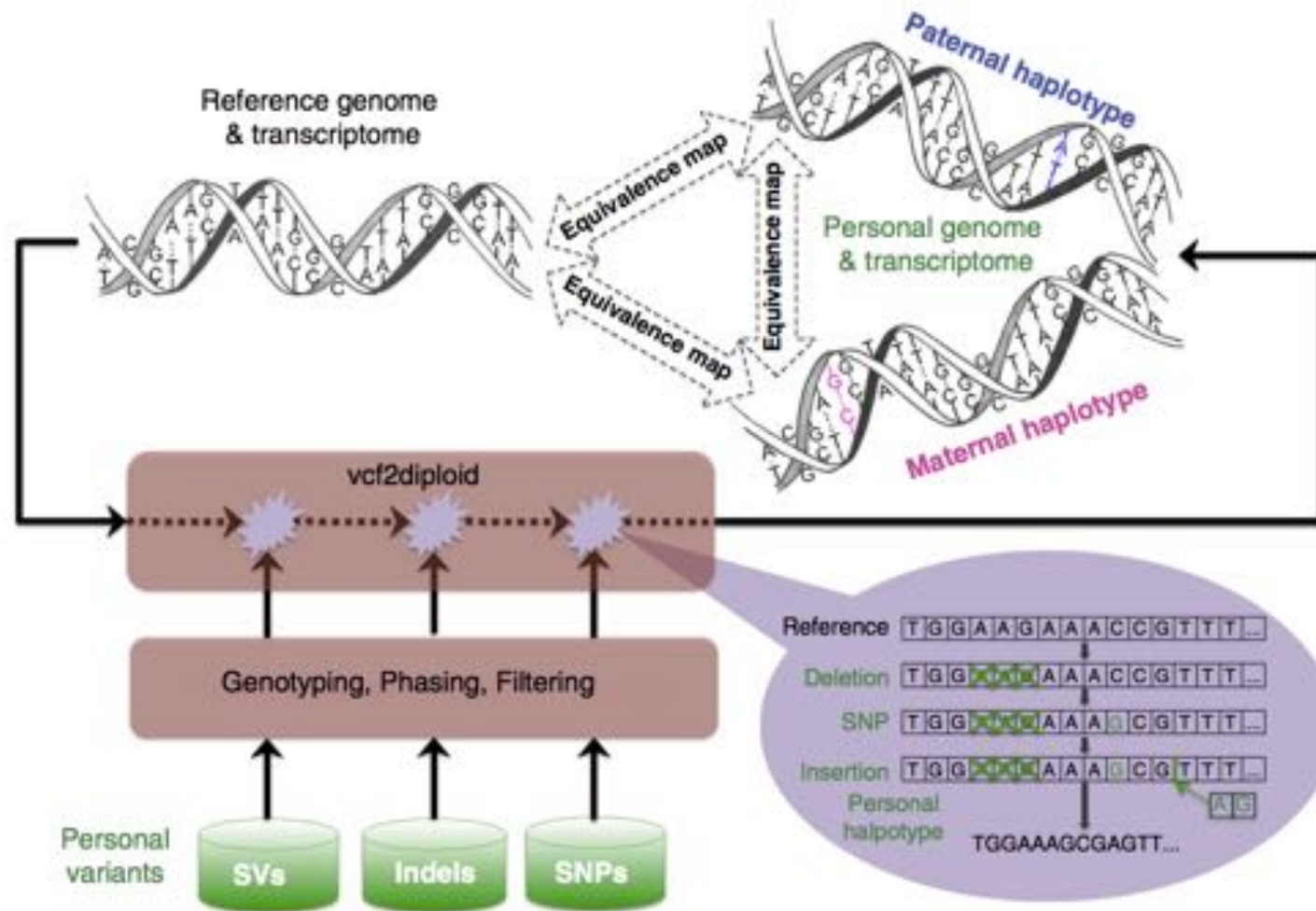
# Hybrid Phasing of Structural Variations

Use the phased short read variants to phase the long reads
The phased long reads allow the SVs to be phased



**Phased Short Read Variants**

**Phased Long Read Variants**

Deletion must be on the orange haplotype!

# Creating a "Perfect" Phased Diploid Genome



(J Rozowsky et al, 2011)

**vcf2diploid inserts phased variants from a VCF file into the reference genome to create a pair of phased chromosome fasta files**

# In pursuit of perfect genome sequencing

1. Why "Perfect"?

2. What is "Perfect"?

3. **How will we achieve it?**
   *Lonnnnnng reads + Looooong mates :-)*
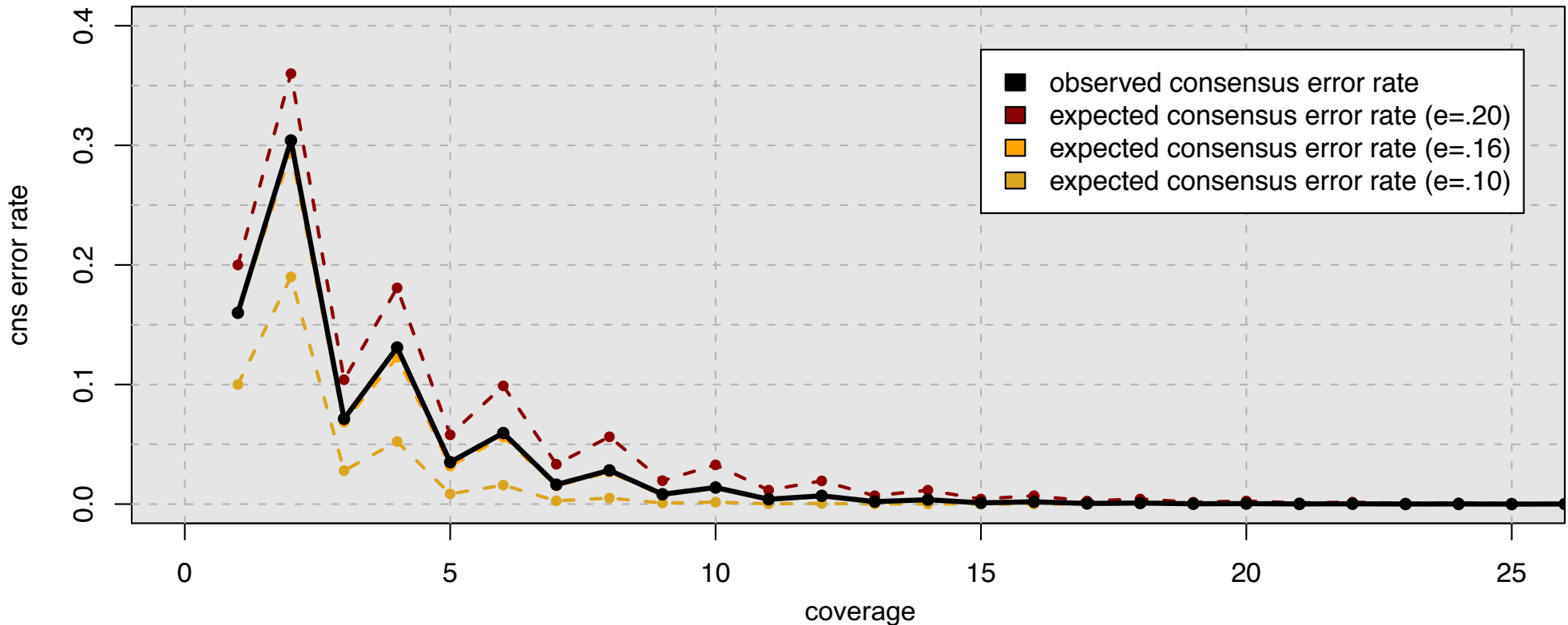
4. When will we achieve it?

# In pursuit of perfect genome sequencing

1. Why "Perfect"?

2. What is "Perfect"?

3. How will we achieve it?

4. **When will we achieve it?**

# Consensus Accuracy and Coverage



# Coverage can overcome random errors

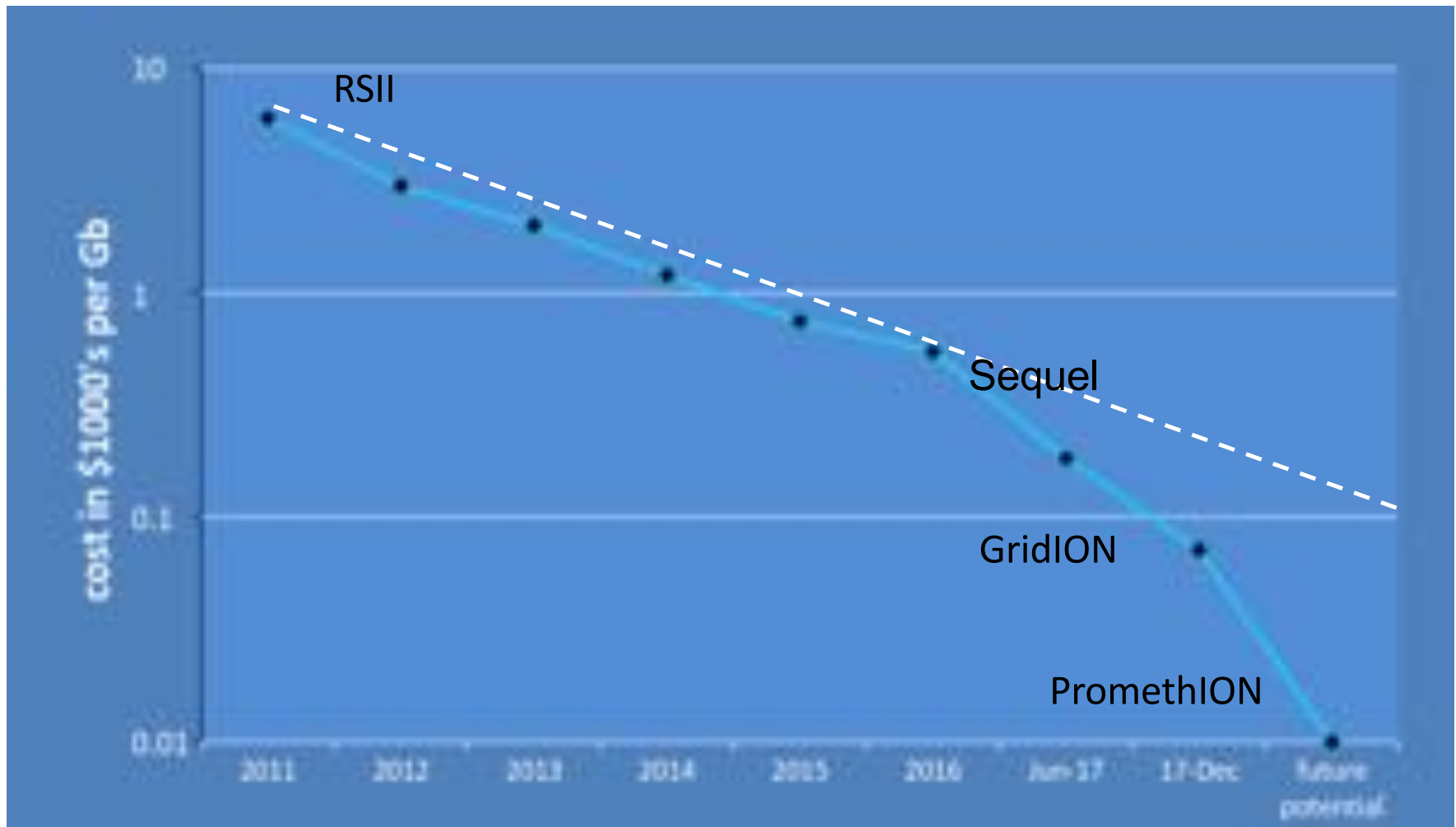- Dashed: error model from binomial sampling
- Solid: observed accuracy

$$CNS\, Error\ =\ \sum_{i=\lceil c/2 \rceil}^{c} \binom{c}{i} (e)^{i} (1-e)^{n-i}$$

*Hybrid error correction and de novo assembly of single-molecule sequencing reads.*
Koren et al (2012) *Nature Biotechnology.* doi:10.1038/nbt.2280

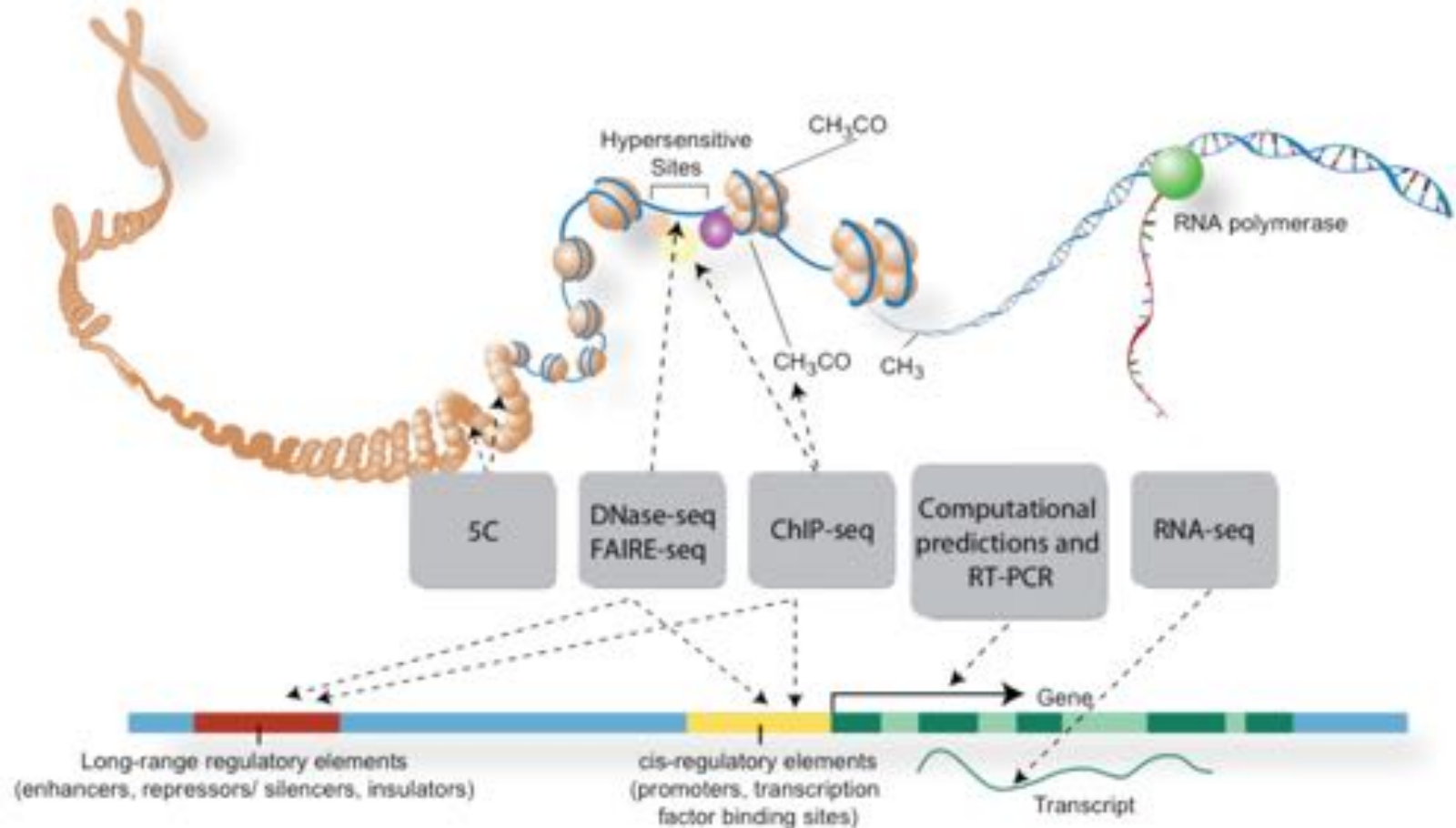# Costs for Long Read Sequencing



Sara Goodwin, CSHL

# "Perfect" Genome Projects



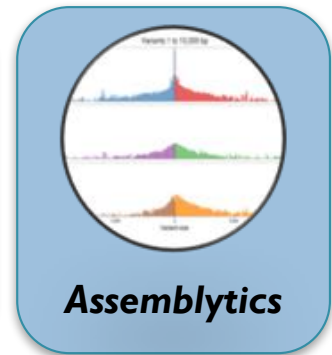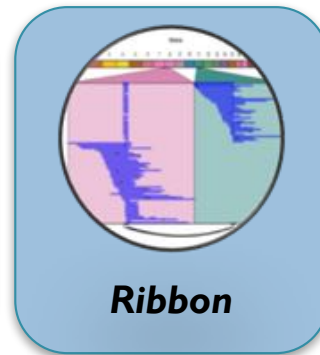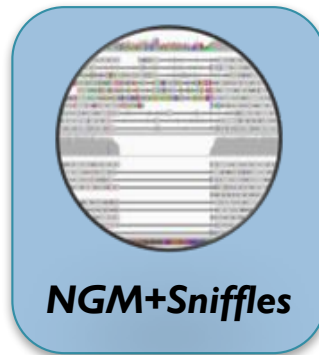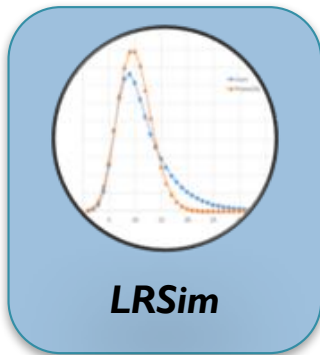| *ENCODE + CancerCODE* | *MaizeCode* | *Tomato Diversity* |
|---|---|---|
| Illumina + PacBio/ONT + 10X RNA-seq, ChipSeq, Hi-C, etc 4 healthy + 10 Organoids | Illumina + PacBio/ONT + 10X RNA-seq, ChipSeq, MNase-seq 2 maize + 2 teosinte | PacBio/ONT + 10x RNA-seq 50 accessions |

# In pursuit of perfect genome sequencing

- *Strive for Perfection: 100% Correct and 100% Complete*
  - The key for perfect genomes is lonnnnnnnnnng reads ☺
  - Expect new insights on the causes of diseases, forces of evolution

- *Multiple sequencing technologies & approaches needed*
  - *PacBio*: Best Resolution of SVs
  - *10X/HIC:* Best Phasing
  - *De novo*: Best Resolution of small SVs
  - *Mapping*: Best resolution of large SVs

- *We have just begun to explore the universe of variants present*
  - Tens of thousands of SVs per person, many megabases of variation
  - Also need to push these ideas into single cell and population scale analysis



**FALCON**  **LRSim**  **SURVIVOR**  **NGM+Sniffles**  **Ribbon**  **Assemblytics**

*http://schatz-lab.org*

# Acknowledgements

# Thank you!

@mike_schatz