# Repetitivní genom, metody NGS za využití dlouhých "readů"

Monika Čechová

@biomonika
biomonika.org

# Interests

- Sex Chromosomes
- Satellite Biology and Heterochromatin
- Long reads and complete genomes
- Early Embryonic Development
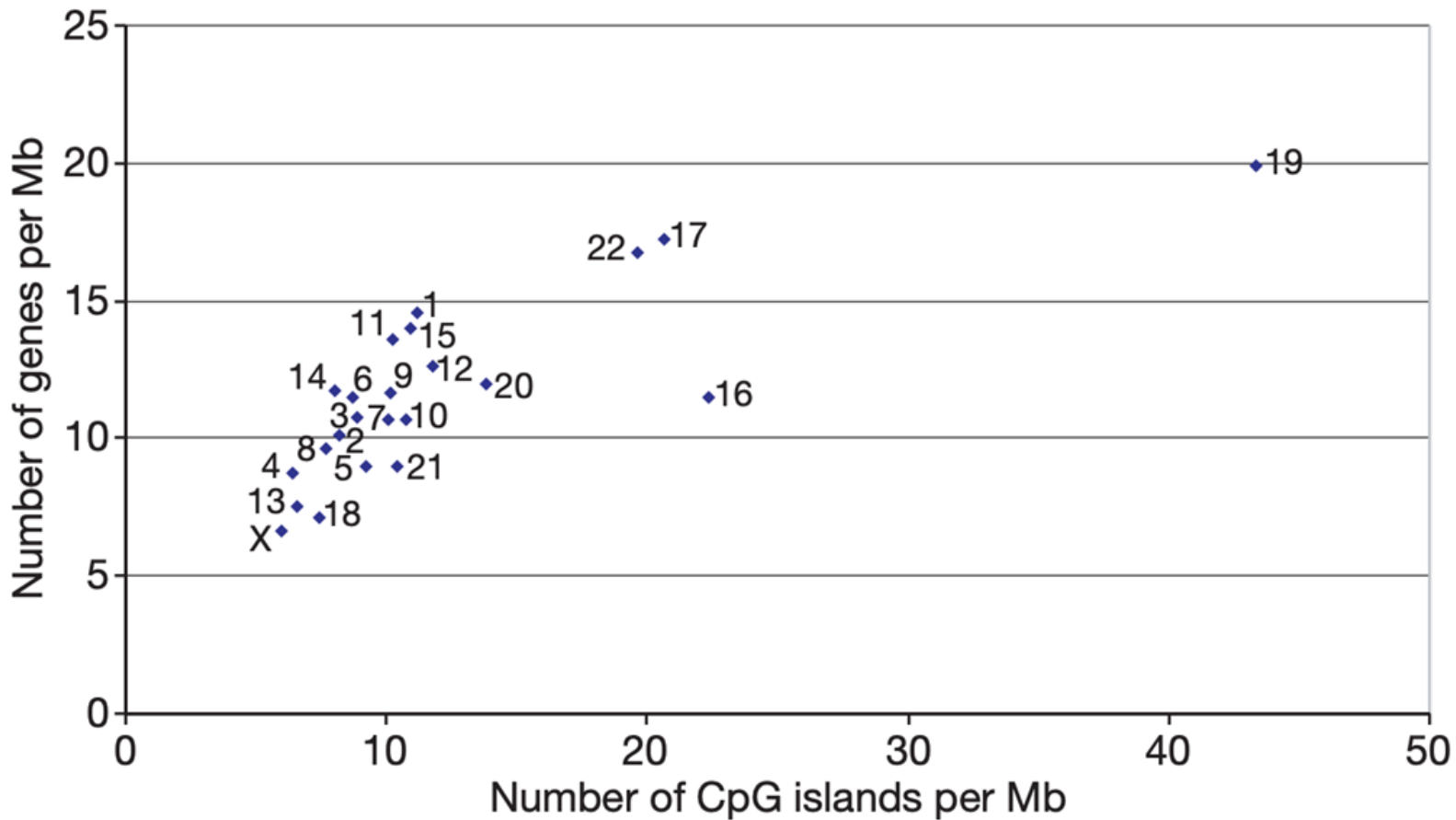- Reproductive Biology

# Education

🎓 PhD Major in Biology, Minor in Statistics, 2020
Penn State, USA

🎓 MS in Bioinformatics, 2013
Masaryk University, Brno

🎓 BS in Applied Informatics, 2011
Masaryk University, Brno

My motivation:
***Answering biological questions using next-generation sequencing data***

- Progress in the human genome assembly and analysis

- Overview of the sequencing technologies

- Long-read applications

- Error correction of long reads

- Structural variation detection: the basics

# Progress in the human genome assembly and analysis

# HUMAN GENOME: WHERE WE STARTED



Hierarchical shotgun sequencing

Genomic DNA

BAC library

Organized mapped large clone contigs

BAC to be sequenced

Shotgun clones

Shotgun sequence
...ACCGTAAATGGGCTGATCATGCTTAAA
                TGATCATGCTTAAACCCTGTGCATCCTACTG...

Assembly  ...ACCGTAAATGGGCTGATCATGCTTAAACCCTGTGCATCCTACTG...

## Classes of interspersed repeat in the human genome



**Figure 17** Almost all transposable elements in mammals fall into one of four classes. See text for details.
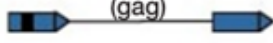
# Homologous sequences as a substrate for a change



Jobling et al., 2008

# Homologous sequences as a substrate for a change



COPY NUMBER CHANGE

INVERSION

GENE CONVERSION

A. Non-allelic homologous recombination

(i) By unequal crossing-over

Duplication

and deletion

NAHR

Hastings et al., 2009

A B C

TEL

Inversion

C B A

TEL

Doug Brutlag, 2011

# Human genome is incomplete



Region containing alternate loci

Region containing fix patches

Region containing novel patches

14

NNNNNNNNNNNNNNN...

https://github.com/broadgsa/gatk/blob/master/doc_archive/dictionary/Reference_Genome_Components

# HUMAN GENOME: WHERE WE ARE HEADING

## Telomere-to-telomere consortium

We have sequenced the CHM13hTERT human cell line with a number of technologies, including 120x coverage of Oxford Nanopore, 70x PacBio CLR, 30x PacBio HiFi, 50x 10X Genomics, as well as BioNano DLS and Arima Genomics HiC. Most raw data is available from this site, with the exception of the PacBio data which was generated by the University of Washington and is available from NCBI SRA.

Human genomic DNA was extracted from the cultured cell line. As the DNA is native, modified bases will be preserved. Nanopore sequencing was performed using Josh Quick's ultra-long read (UL) protocol.

The X chromosome was selected for manual assembly, and was initially broken at three locations: the centromere (artificially collapsed in the assembly), a large segmental duplication (DMRTC1B, 120 kb), and a second segmental duplication with a paralogue on chromosome 2 (134 kb). Gaps in the GRCh38 reference (black) and known segmental duplications (red; paralogous to Y, pink) are annotated. Repeats larger than 100 kb are named with the expected size (kb) (blue, tandem repeats; red, segmental duplications).

# Overview of the sequencing technologies

# Outline



https://www.gatc-biotech.com/fileadmin/_processed_/
csm_Sanger_sequencing_illustration_small_898122494c.jpg

# Illumina

- Introduction (5 mins)
  - https://www.youtube.com/watch?v=fCd6B5HRaZ8

# Applications

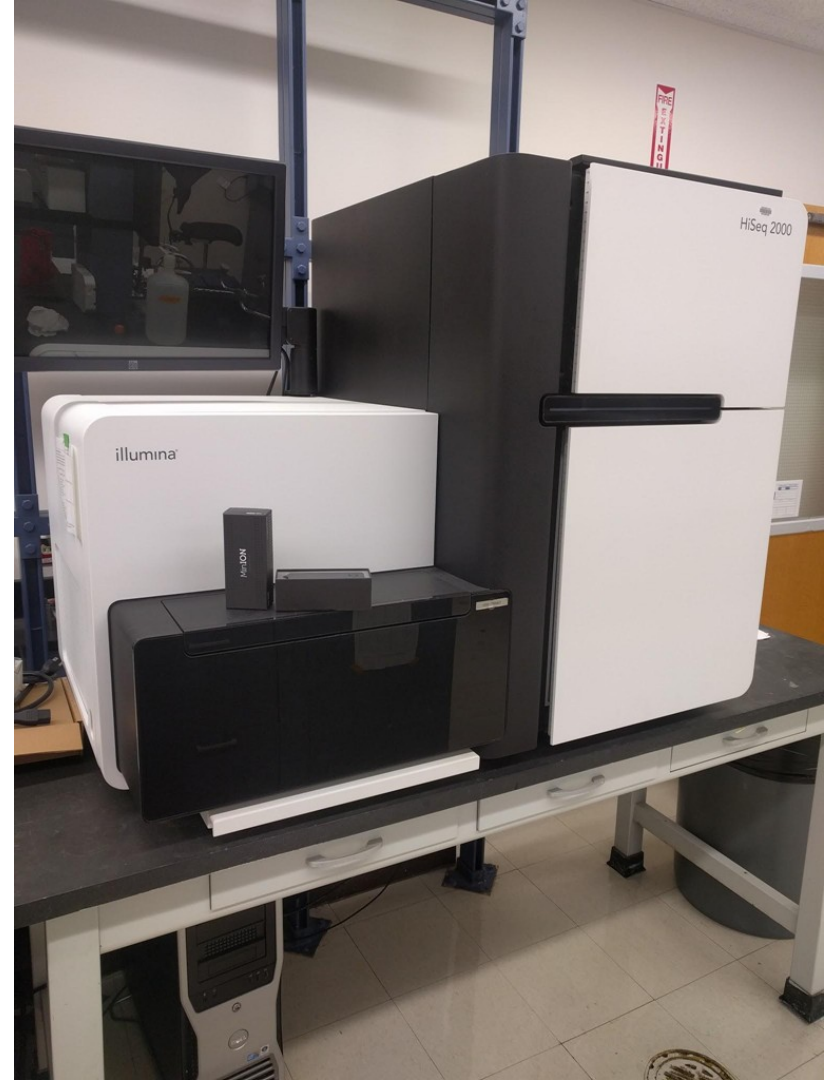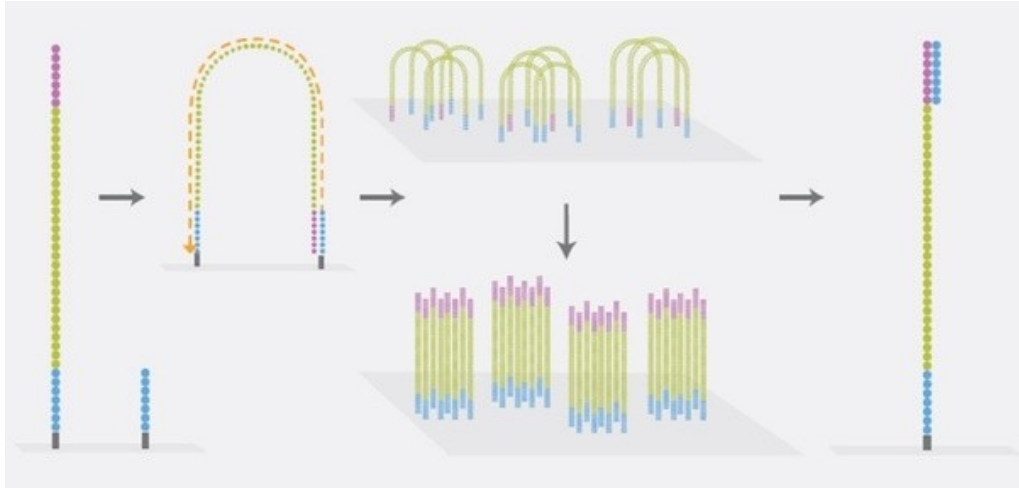## Large-scale whole-genome sequencing of the Icelandic population

Daniel F Gudbjartsson ✉, Hannes Helgason […] Kari Stefansson ✉

### Abstract

Here we describe the insights gained from sequencing the whole genomes of 2,636 Icelanders to a median depth of 20×. We found 20 million SNPs and 1.5 million insertions-deletions (indels). We describe the density and frequency spectra of sequence variants in relation to their functional annotation, gene position, pathway and conservation score. We demonstrate an excess of homozygosity and rare protein-coding variants in Iceland. We imputed these variants into 104,220 individuals down to a minor allele frequency of 0.1% and found a

Sequencing versus genotyping?

# Applications



Normal nuchal translucency | Increased nuchal translucency

http://s19.postimage.org/nx6ifenv7/2012_10_11_215737.png

## NIPT Is Noninvasive to the Mother and Baby

NIPT analyzes cell free DNA from a maternal blood sample (mixture of fetal and maternal DNA) to screen for common chromosomal conditions including trisomy 21 (Down syndrome), trisomy 18 (Edwards syndrome), and trisomy 13 (Patau syndrome).

illumina®

# PacBio

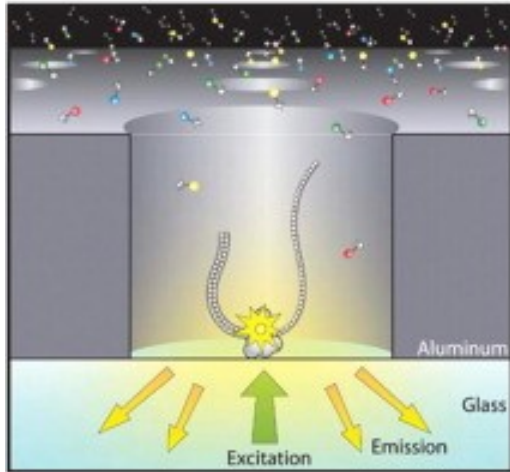- Introduction to SMRT sequencing (2 mins)
  - https://www.youtube.com/watch?NHCJ8PtYCFc
- Advanced overview of the technology
  - https://vimeo.com/121267426 (9 mins)
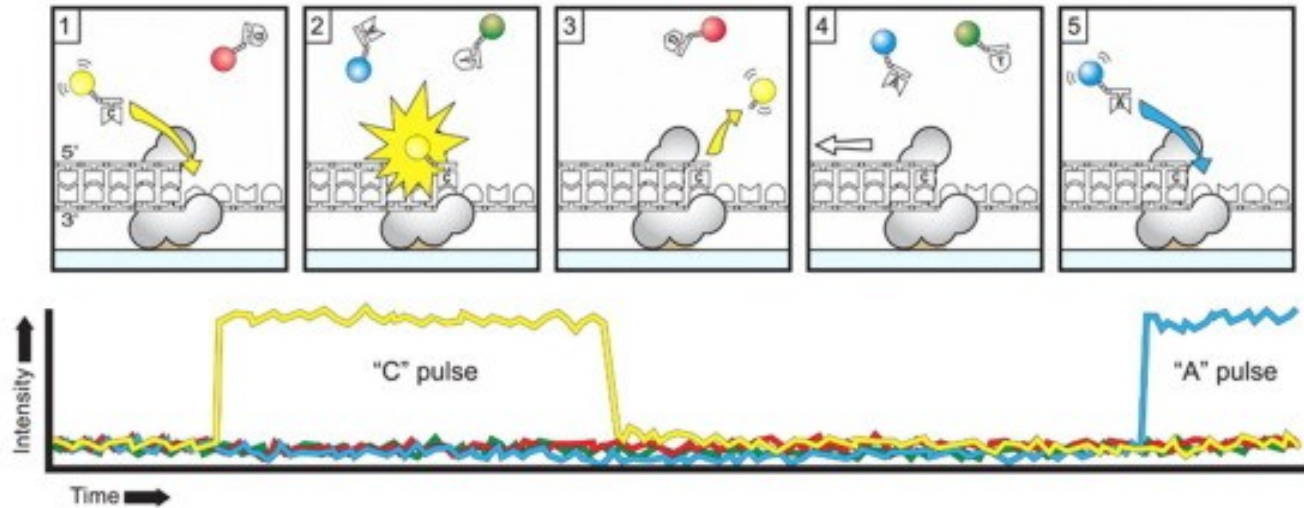


https://www.pacb.com/wp-content/uploads/sequel-2Bimage-300x295.jpg

# PacBio

# Nanopore

# Nanopore

- Introduction (3 mins)
  - https://www.youtube.com/watch?v=GUb1TZvMWsw
- Advanced slides on the technology
  - https://github.com/lmmx/talk-transcripts/blob/master/Nanopore/NoThanksIveAlreadyGotOne.md

https://twitter.com/MakovaLab

# Nanopore

- Introduction (3 mins)
  - https://www.youtube.com/watch?v=GUb1TZvMWsw
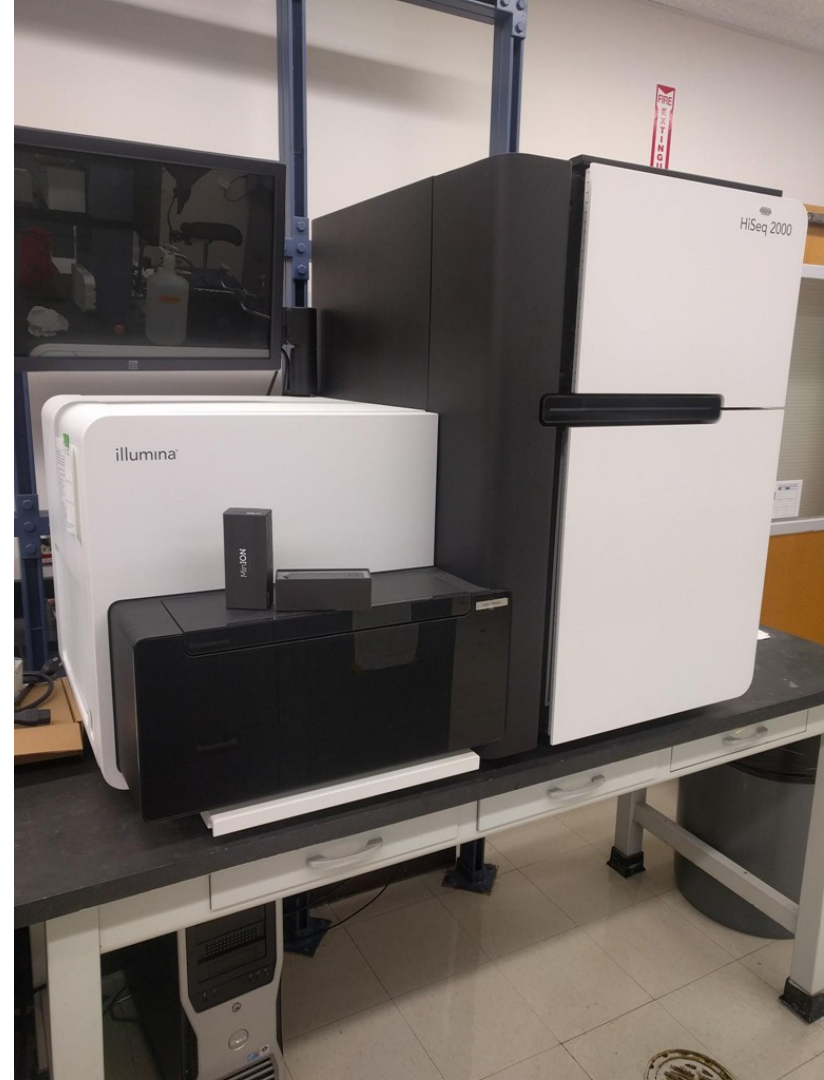- Advanced slides on the technology
  - https://github.com/lmmx/talk-transcripts/blob/master/Nanopore/NoThanksIveAlreadyGotOne.md

https://twitter.com/MakovaLab

# Real-time surveillance

## Real time genomic surveillance of Ebola outbreak 2014-2015

05 Jun 2015

The current Ebola outbreak in West Africa is the largest ever recorded, with over 26,500 cases reported resulting in an estimated 11,000 deaths. Yet genomic surveillance of this outbreak has been patchy, hampered by understandable but vexing logistical, social, political and technical obstacles in securing and transporting samples for processing.

We wanted to help address the gaps in our knowledge of viral evolution and to generate data for epidemiological use. So, in April, Josh Quick from my group went to Conakry, Guinea to establish proof-of-principle for portable nanopore sequencing. This was the most practical way we could rapidly establish a local sequencing lab in order to generate real-time information.

## Rapid draft sequencing and real-time nanopore sequencing in a hospital outbreak of *Salmonella*

Joshua Quick [†], Philip Ashton [†], Szymon Calus, Carole Chatt, Savita Gossain, Jeremy Hawker, Satheesh Nair, Keith Neal, Kathy Nye, Tansy Peters, Elizabeth De Pinna, Esther Robinson, Keith Struthers, Mark Webber, Andrew Catto, Timothy J. Dallman, Peter Hawkey ✉ and Nicholas J. Loman ✉

[†]Contributed equally

# Sequencing the station: Investigation aims to identify unknown microbes in space



Fig. 1 Astronaut Kate Rubins on the ISS

https://phys.org/news/2017-04-sequencing-station-aims-unknown-microbes.html



a) First of four datasets generated on Earth as ground controls

b) First of four datasets generated on ISS
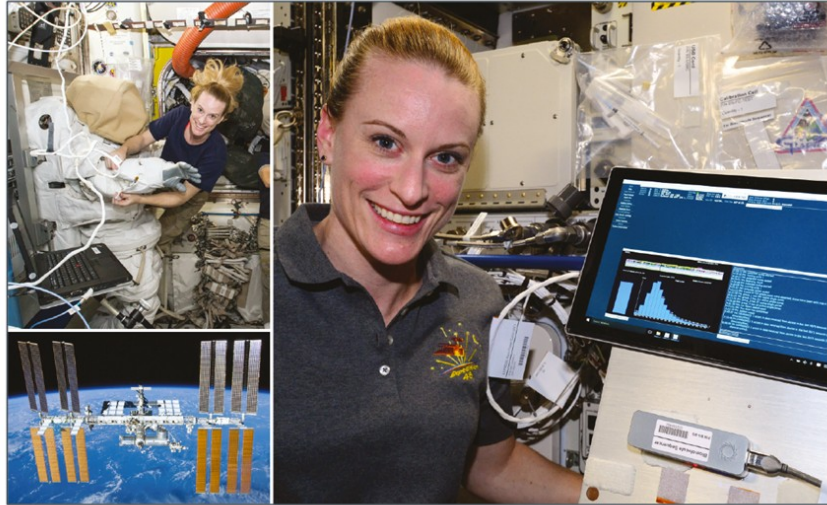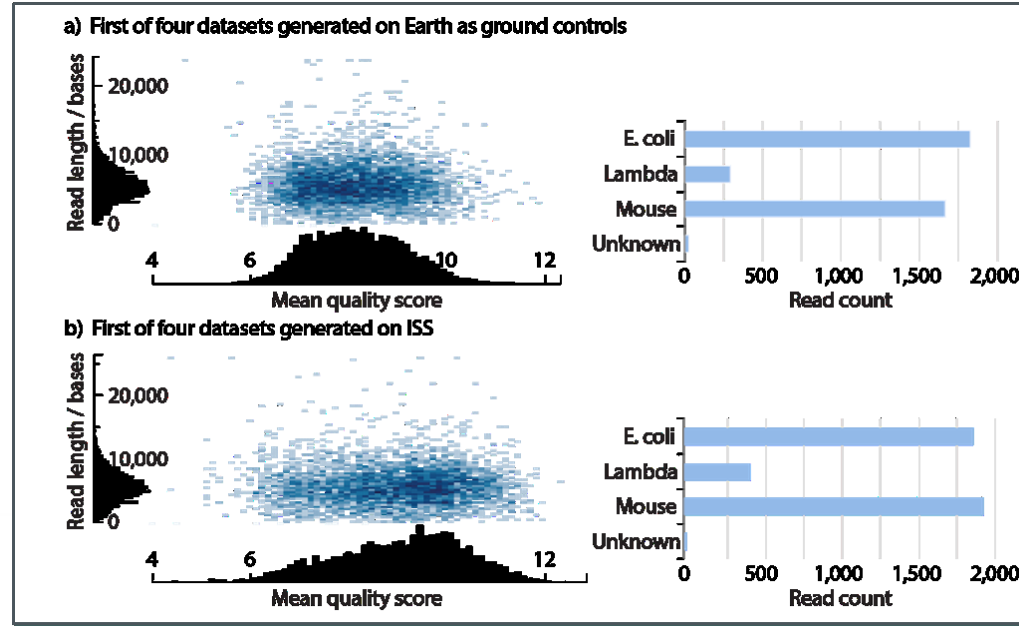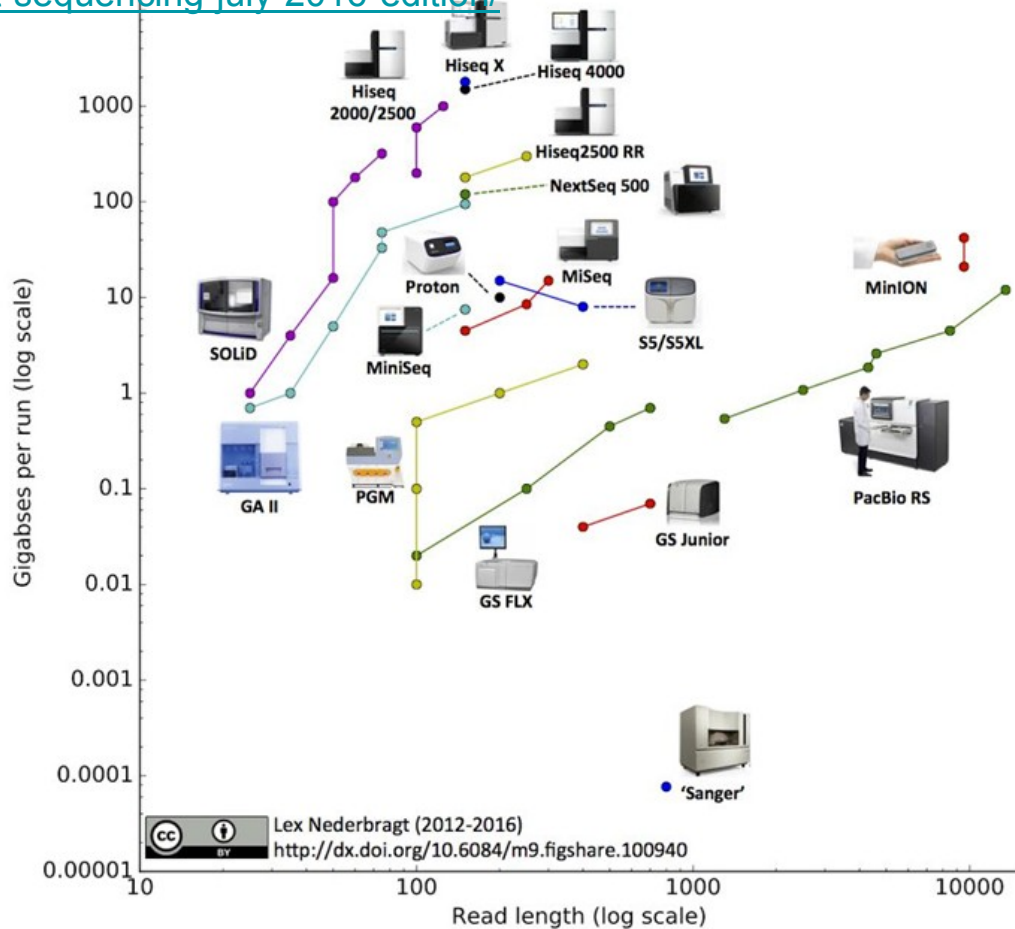
Fig. 2 Analysis workflow showing read quality for runs a) on Earth b) on the ISS

https://nanoporetech.com/resource-centre/posters/dna-sequencing-microgravity-international-space-station-iss-using-minion

Lex Nederbragt (2012-2016)
http://dx.doi.org/10.6084/m9.figshare.100940

# Long-read sequencing

- Sample quality
- Library preparation (size selection, repair)
- Protocols
- Sequencing throughput
- Multiplexing

Budgeting is important: get yourself familiar with the cost of reagents and sequencing!
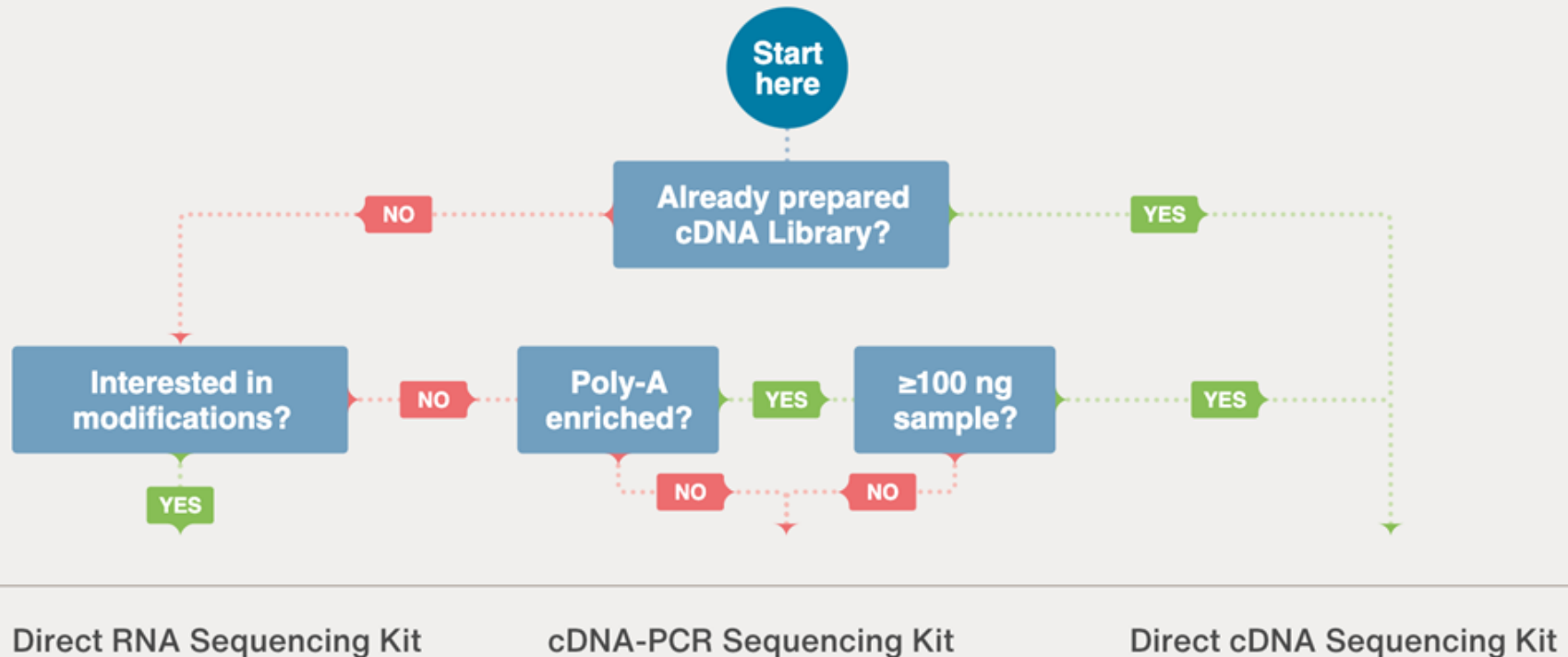
**There's always a trade-off.**
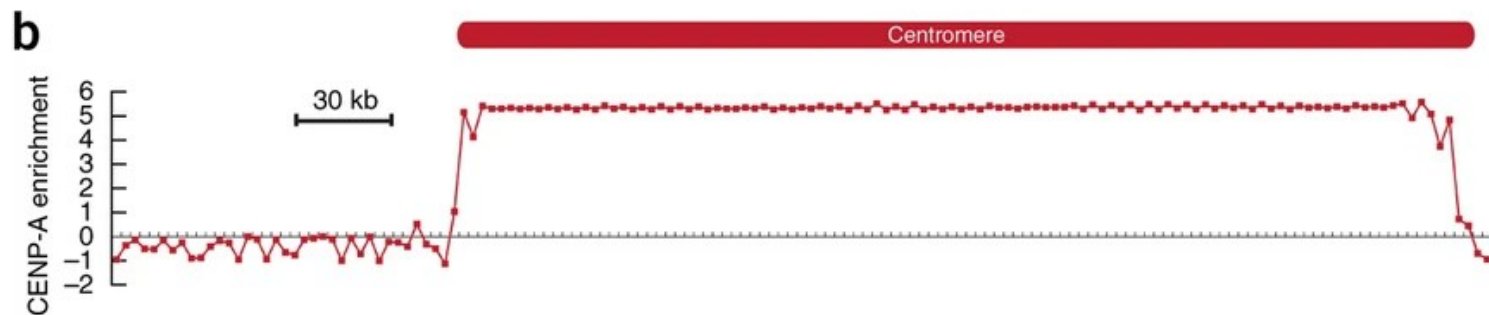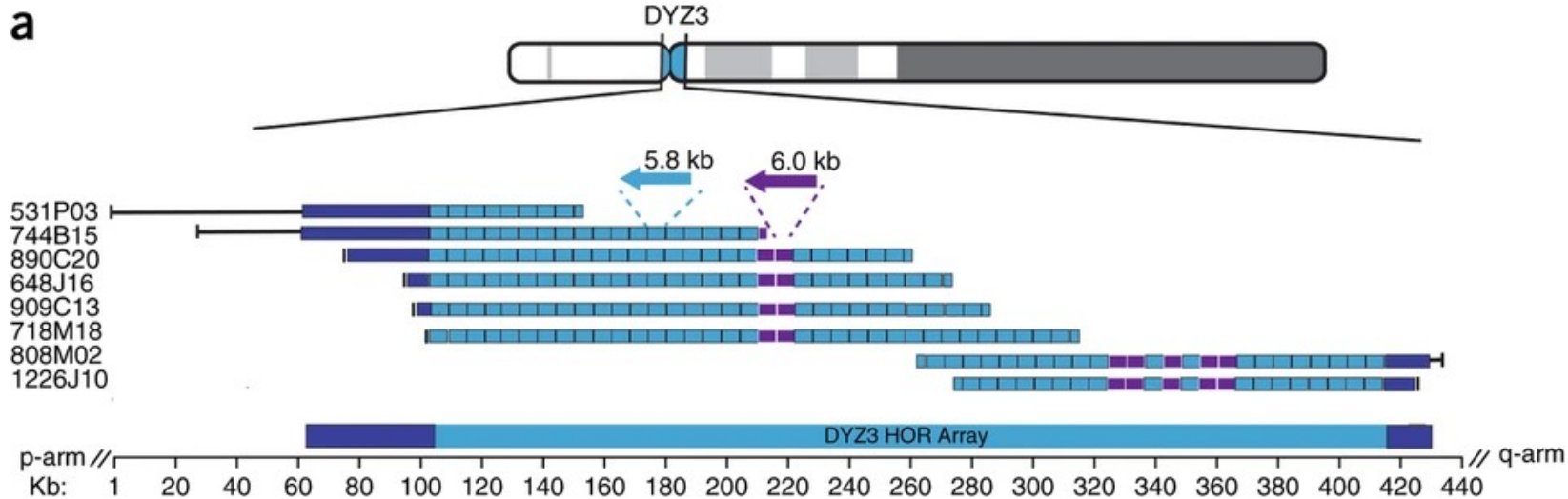
# Long-read applications

# Long reads

1. **Transcriptome (isoforms)**
2. **Assembly (PacBio HiFi + ultra-long Nanopore reads)**
3. **Epigenetic modifications**

# 1. RNA sequencing with Nanopore

# 2. Linear assembly of the RP11 Y centromere

# 3. Nanopore sequencing meets epigenetics

Figure 1 | DNA methylation can be read out directly by nanopore sequencing. Single-stranded DNA alters ionic current in a sequence-dependent manner as it passes through a pore (left), with methylated bases highlighted with red and blue tags. The raw current signal (right) indicates small changes due to methylation that a new set of algorithms can robustly interpret.

# Error correction of long reads

# Error profiles

**Table S1. Observed error rates in PacBio and Nanopore alignments.** m=matches, mm=mismatches, io=insertion open, ix=insertion extend, do=deletion open, dx=deletion extend. Overall error rates are measured as 1 - m/(m+mm+io+ix+do+dx).
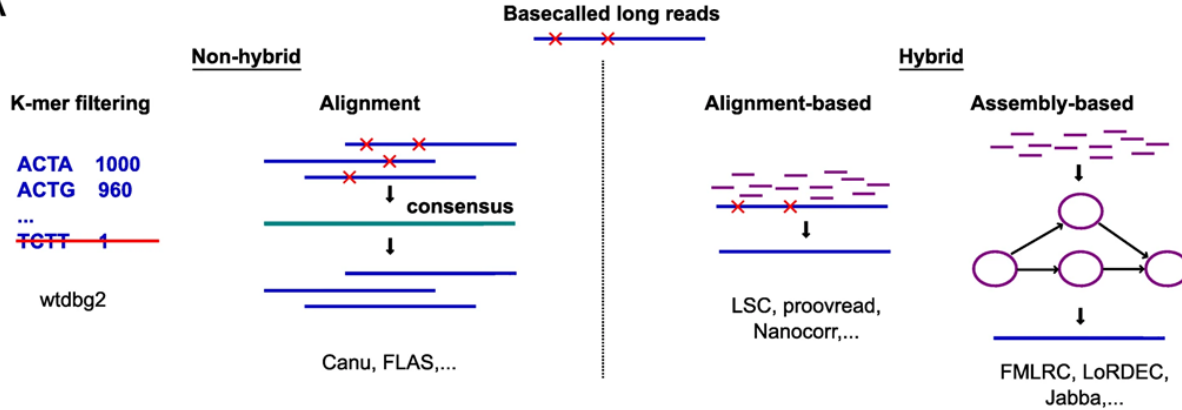
| Technology | Error Rate | m | mm | io | ix | do | dx |
|---|---|---|---|---|---|---|---|
| PacBio | 14.90% | 203521495440 | 4039815573 | 15516830204 | 5868396512 | 8622457071 | 1709970892 |
| | | 85.06% | 1.69% | 6.48% | 2.45% | 3.60% | 0.71% |
| ON | 16.10% | 17603088986 | 965110092 | 492506288 | 299904008 | 702350550 | 919627709 |
| | | 83.89% | 4.60% | 2.35% | 1.43% | 3.35% | 4.38% |

RS Harris, **M Cechova**, KD Makova. *Noise-Cancelling Repeat Finder: Uncovering tandem repeats in error-prone long-read sequencing data. BIOINFORMATICS*, 2019.

# Error correction

- Error correction with short reads
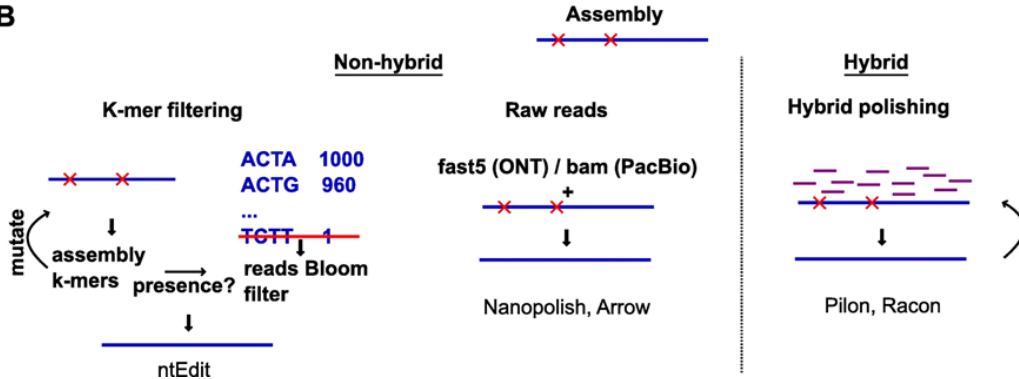- Error correction with long reads
- Assembly merging

# Error correction


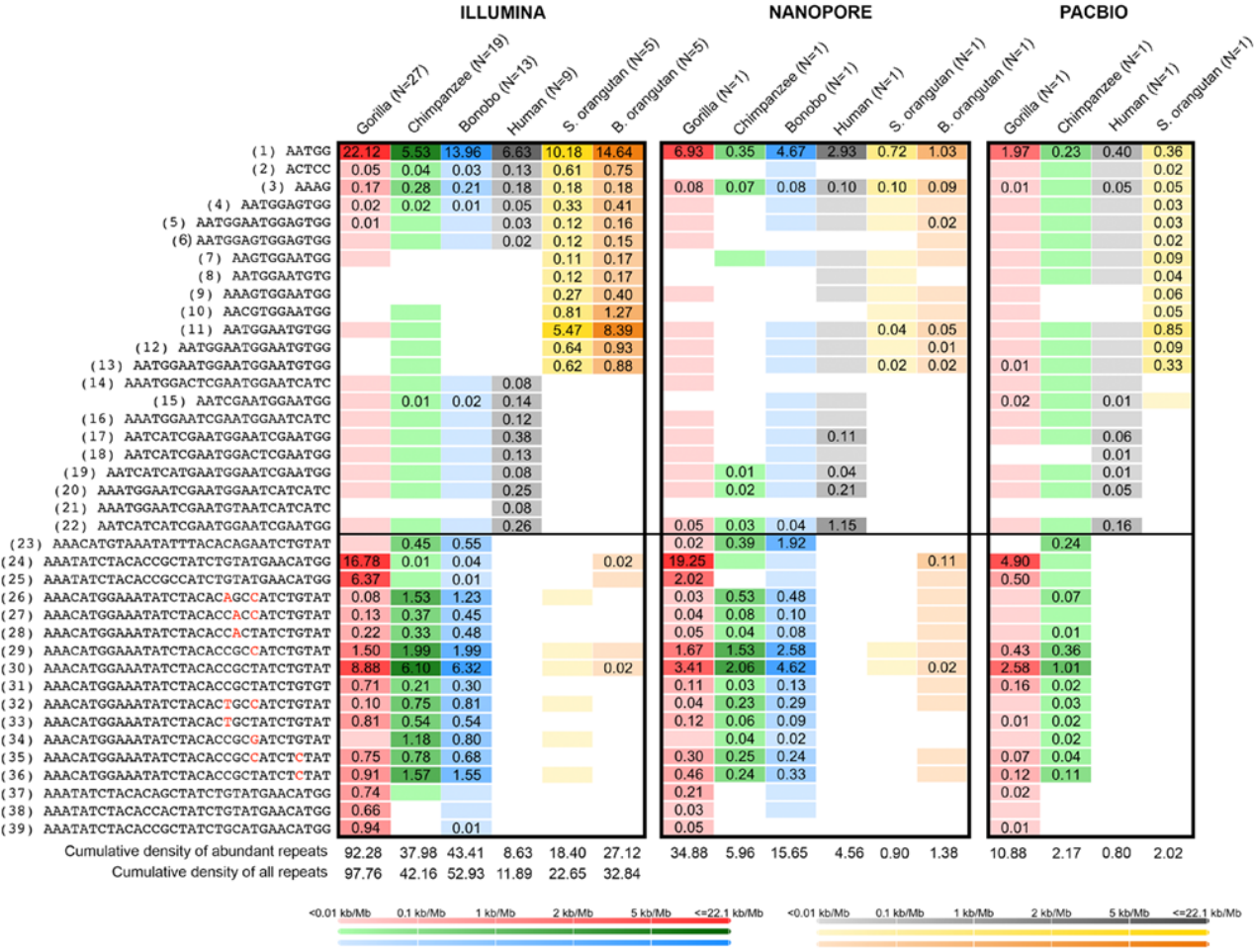
- Error correction with short reads
- Error correction with long reads
- Assembly merging

Amarasinghe, S.L., Su, S., Dong, X. *et al*. Opportunities and challenges in long-read sequencing data analysis. *Genome Biol* 21, 30 (2020). https://doi.org/10.1186/s13059-020-1935-5
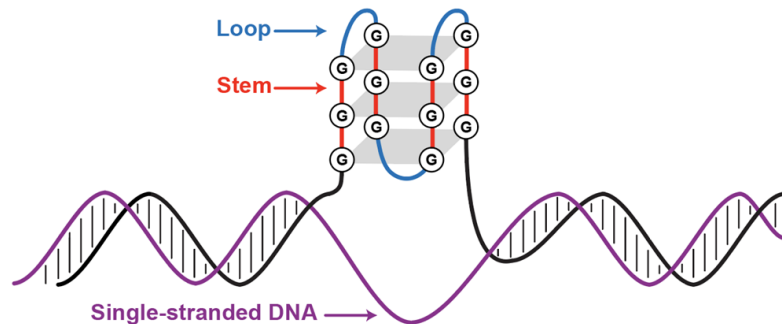
We demonstrated that orthogonal technologies (such as Illumina, Nanopore, and PacBio) are generally concordant in distinguishing between highly and lowly abundant repeated motifs.
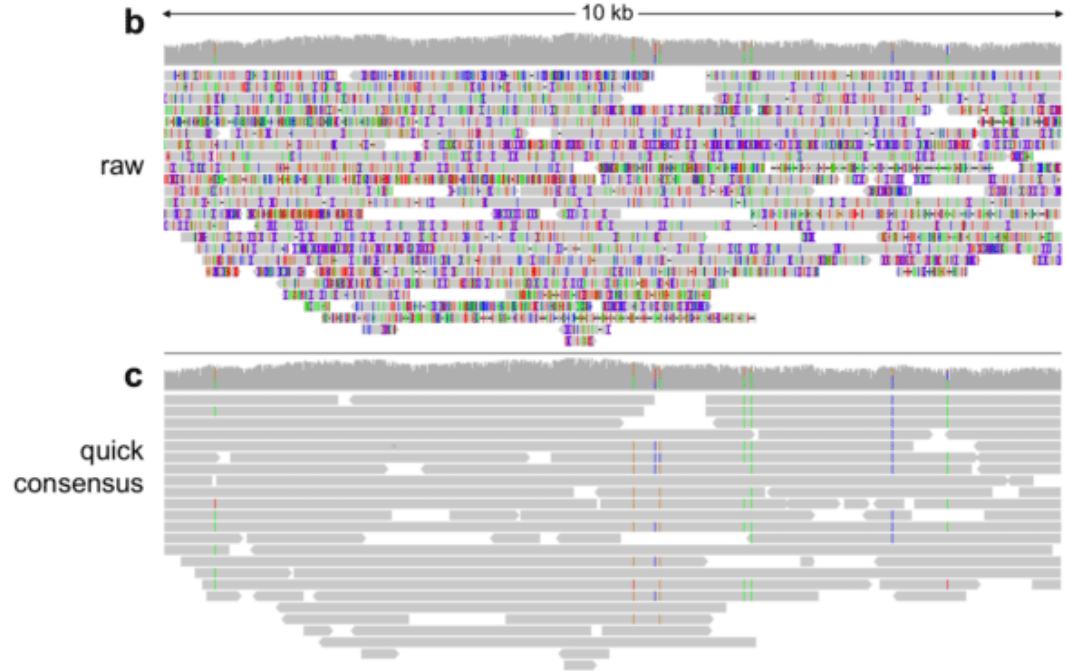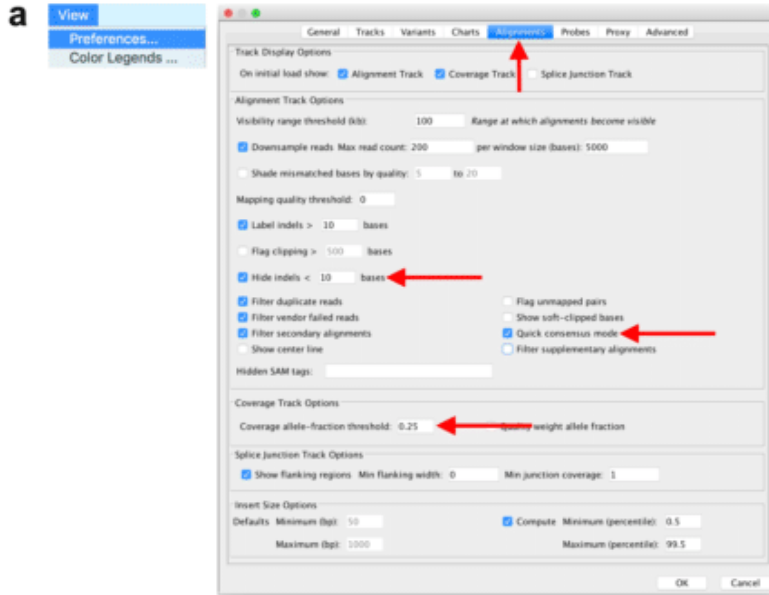


Cechova et al., 2019
Harris et al., 2019

# Long-read sequencing technology indicates genome-wide effects of non-B DNA on polymerization speed and error rate

Wilfried M. Guiblet[1,8], Marzia A. Cremona[2,8], Monika Cechova[3], Robert S. Harris[3],
Iva Kejnovská[4], Eduard Kejnovsky[5], Kristin Eckert[6], Francesca Chiaromonte[2,7] and
Kateryna D. Makova[3]

Loop

Stem

Single-stranded DNA



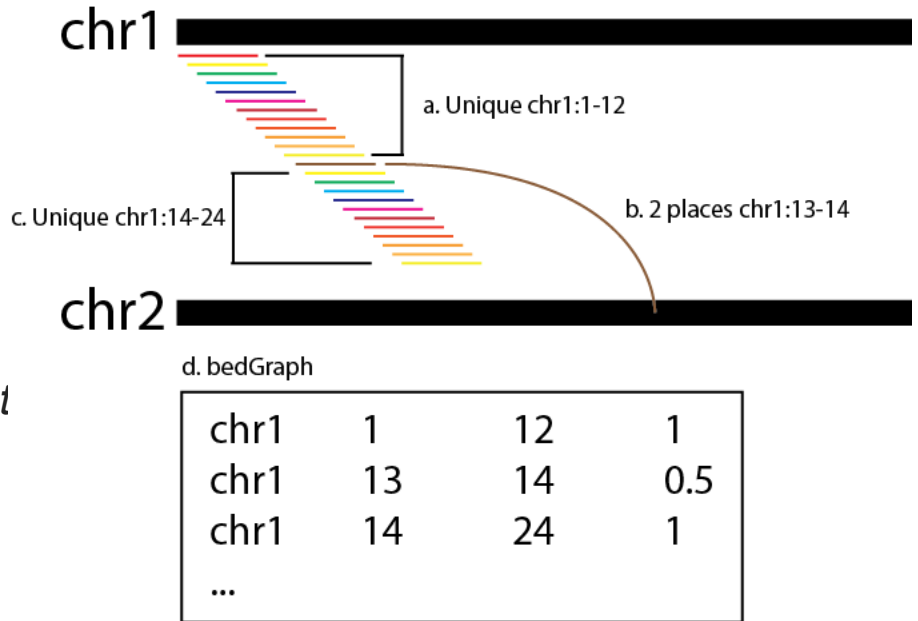SMRT™bell

Consensus Sequence

# The visualization of sequencing errors

# Mappability profiles

*Alignability* – *These tracks provide a measure of how often the sequence found at the particular location will align within the whole genome. Unlike measures of uniqueness, alignability will tolerate up to 2 mismatches. These tracks are in the form of signals ranging from 0 to 1 and have several configuration options.*
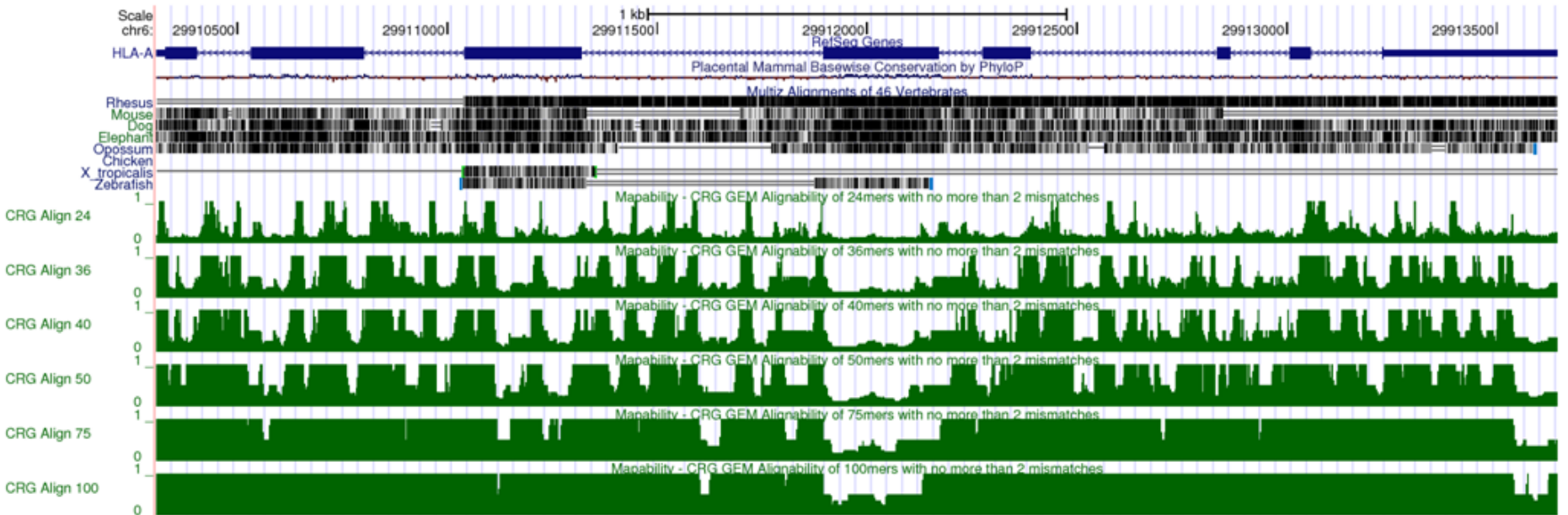
*Uniqueness* – *These tracks are a direct measure of sequence uniqueness throughout the reference genome. These tracks are in the form of signals ranging from 0 to 1 and have several configuration options*



chr1

a. Unique chr1:1-12

b. 2 places chr1:13-14

c. Unique chr1:14-24

chr2

d. bedGraph

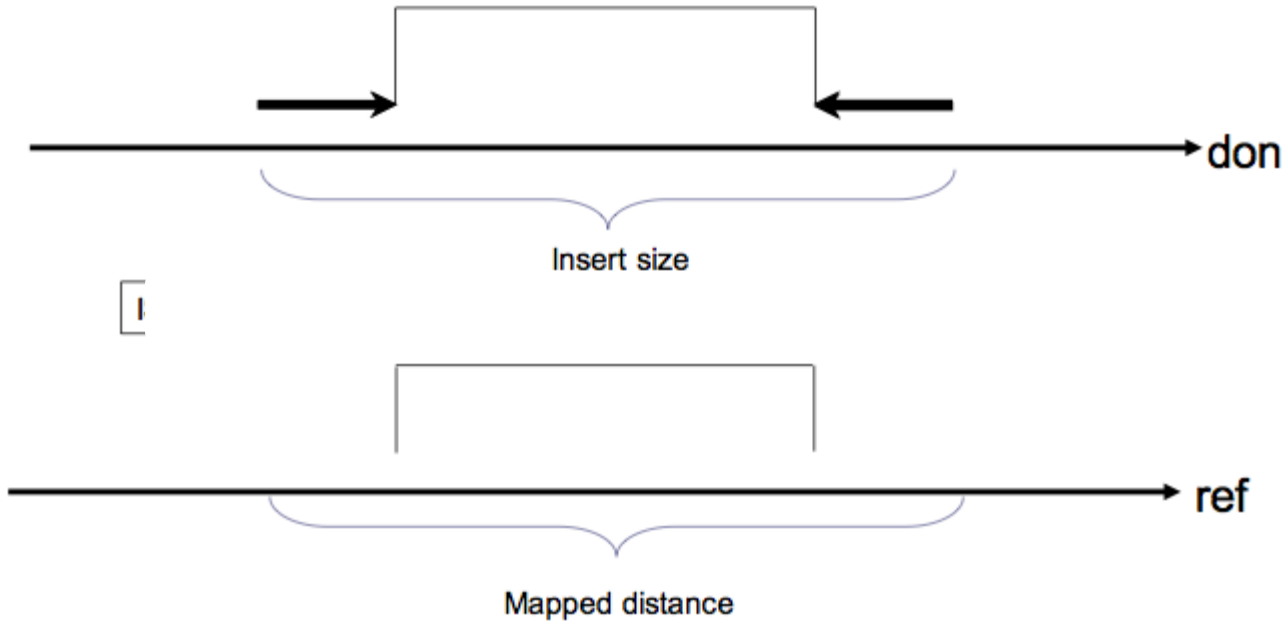| chr1 | 1  | 12 | 1   |
|------|----|----|-----|
| chr1 | 13 | 14 | 0.5 |
| chr1 | 14 | 24 | 1   |
| ...  |    |    |     |

# Mappability profiles

Influence of paralogous genes on the mappability scores: the example of the HLA-A gene.

The HLA-A gene is part of the Major Histocompatibility Complex (MHC) involving a large gene family with numerous paralogs. This screenshot of the UCSC genome browser (with the six mappability tracks in green) illustrates the low uniqueness of the HLA-A gene (especially, its exon 4) which could render its targeting by RNASeq difficult (if only uniquely mapping reads are considered).
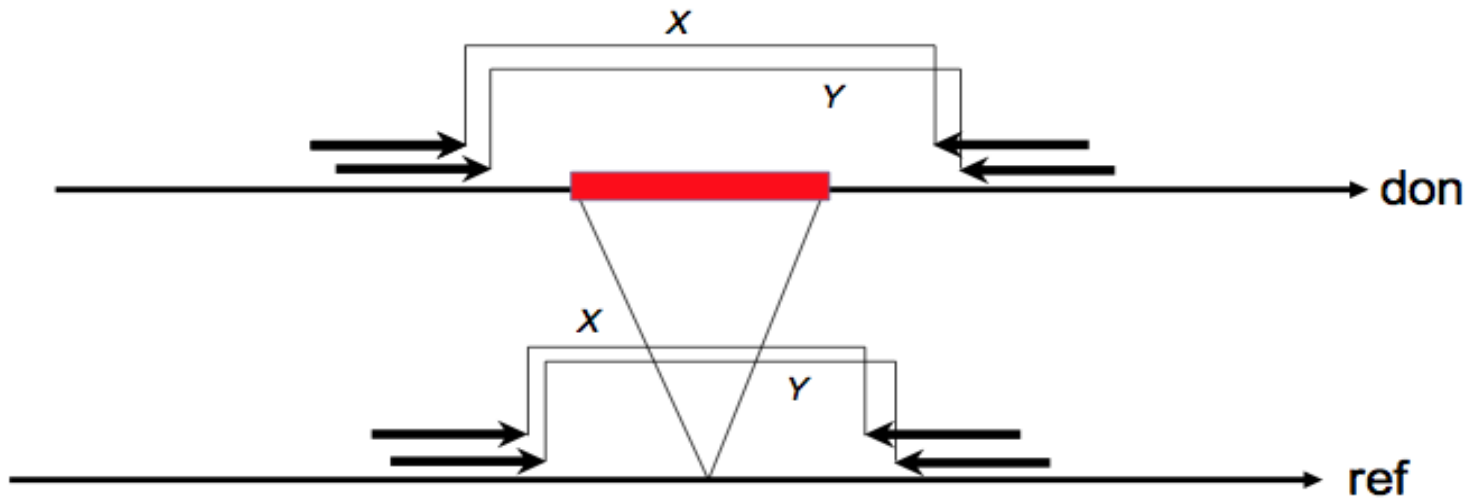


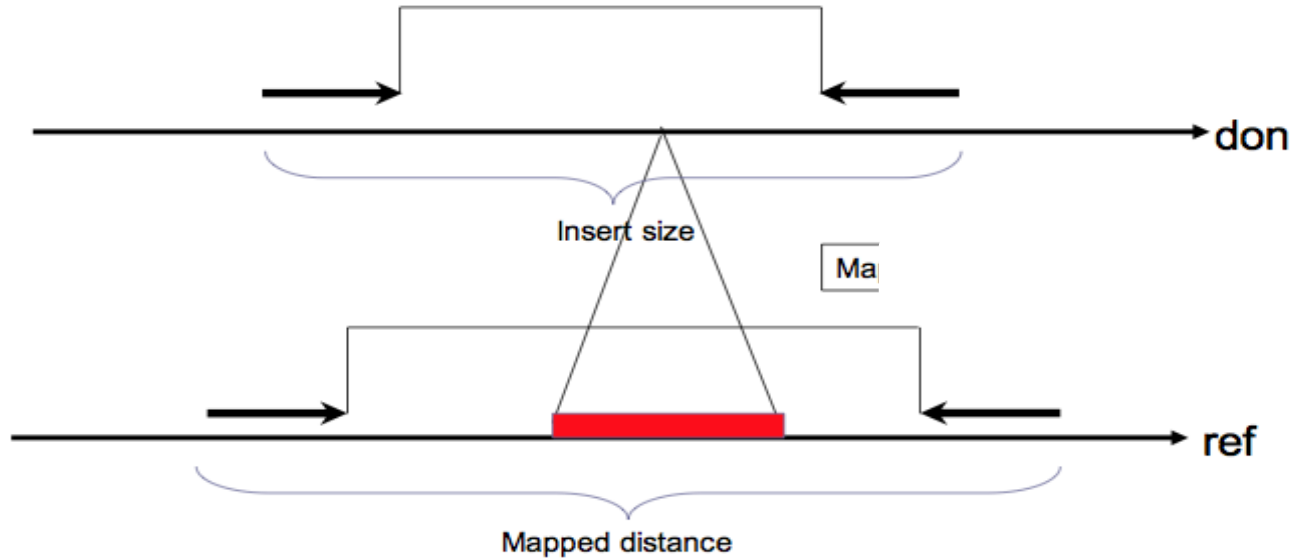doi: https://doi.org/10.1371/journal.pone.0030377.g009

# Structural variation detection: the basics

# Matepair signature: no SVs

# Insertion: consistency



adapted from Prof. Paul Medvedev, Algorithms & Data Structures in Bioinformatics

# Deletion: signature



adapted from Prof. Paul Medvedev, Algorithms & Data Structures in Bioinformatics

# Inversion: signature



Prof. Paul Medvedev, Algorithms & Data Structures in Bioinformatics

# SV summary

| Type | Mapped Distance | Orientation |
|------|-----------------|-------------|
| Insertion | too big | correct |
| Deletion | too small | correct |
| Inversion | * | → ⌐ → |
| Tandem duplication | * | ← → |
| Interchromosomal | different chromosomes | N/A |

Prof. Paul Medvedev, Algorithms & Data Structures in Bioinformatics

# Depth-of-coverage



Depth-of-coverage can
help detect SVs

adapted from Prof. Paul Medvedev, Algorithms & Data Structures in Bioinformatics

# IGV



- view as pairs

http://www.broadinstitute.org/igv

# IGV

- split screen view



http://www.broadinstitute.org/igv

# Hands-on session (not graded)

- run TRF and NCRF and compare the results for the GGAAT repeat

- visualize the alpha satellites in the centromeric portion of the human Y chromosome

- solve the Rosalind problem "Error correction in reads"

# Thank you for your attention