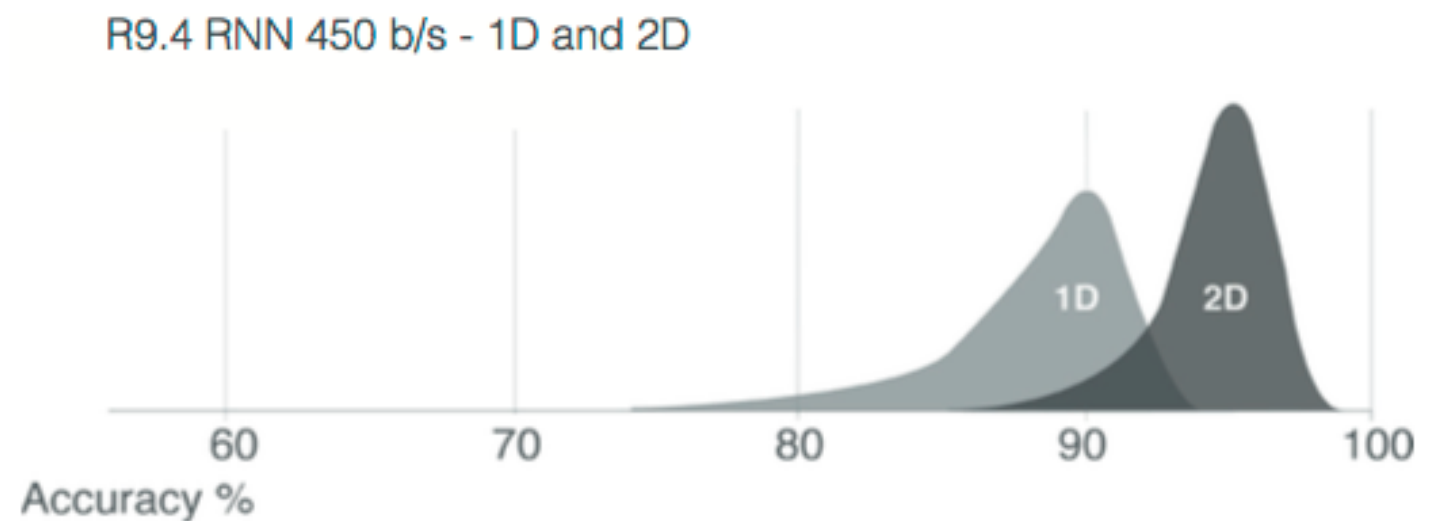# Musing research activities with long read sequence data

Yosuke Tanigawa (ytanigaw@stanford.edu)
M. Rivas Lab (rotation student) | Biomedical Informatics Ph.D. Program

# Oxford Nanopore MinIon

- Handy size DNA sequencer

- Long read sequencer (8kb)

- 3 protocols: 1D, 2D, Rapid

- High error rate
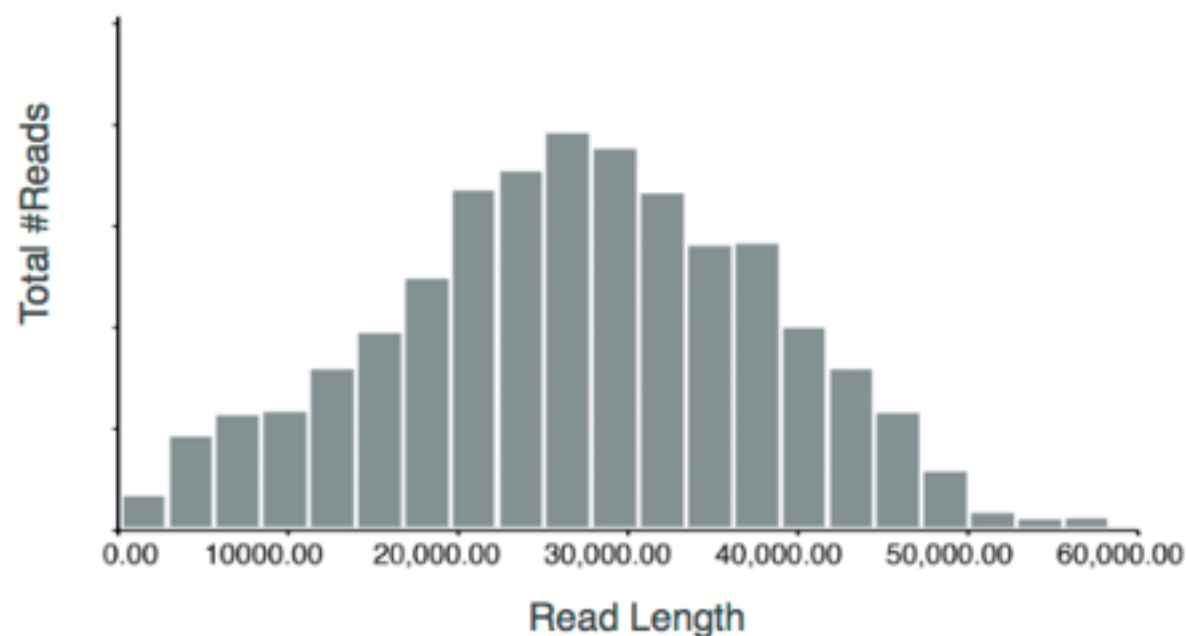
  - 92% (1D), 96% (2D)

- R9.4 chemistry



R9.4 RNN 450 b/s - 1D and 2D

Accuracy %

https://nanoporetech.com/products/minion
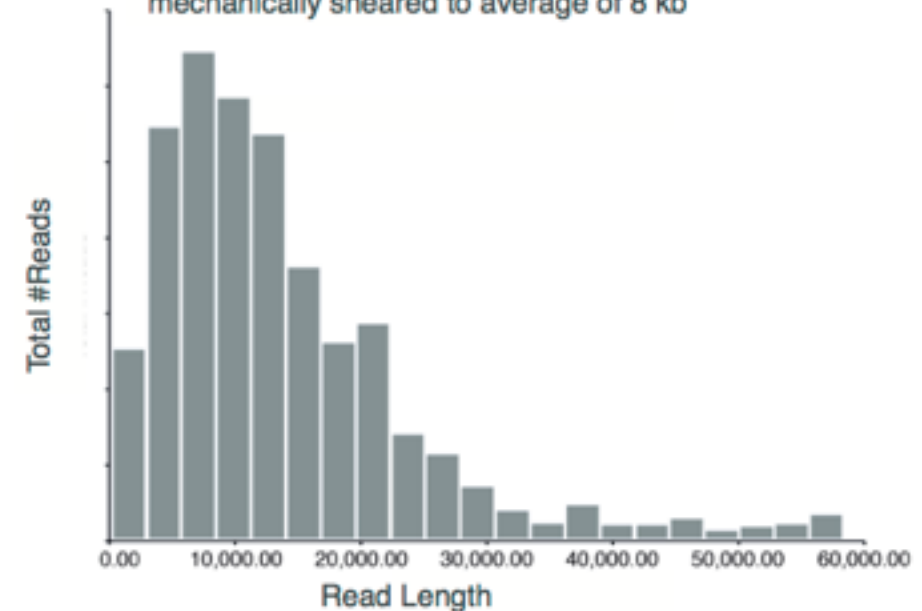
# Read Length Distribution

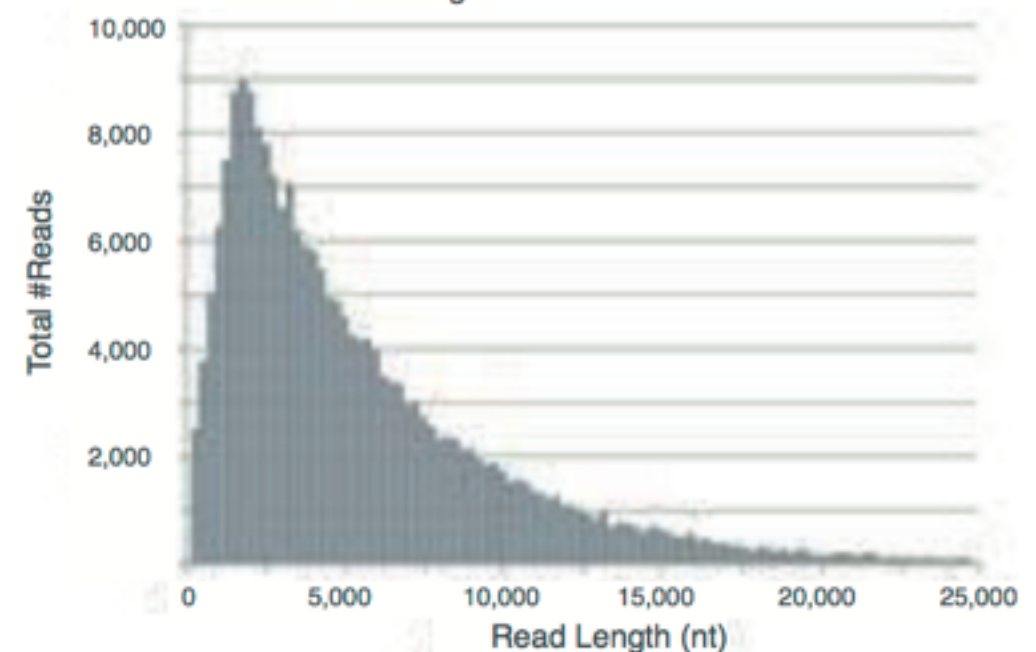Mode of the distributions
- 1D:     8kb
- 2D:   28kb
- rapid: 3kb

1D reads – typical profile of lambda DNA
mechanically sheared to average of 8 kb
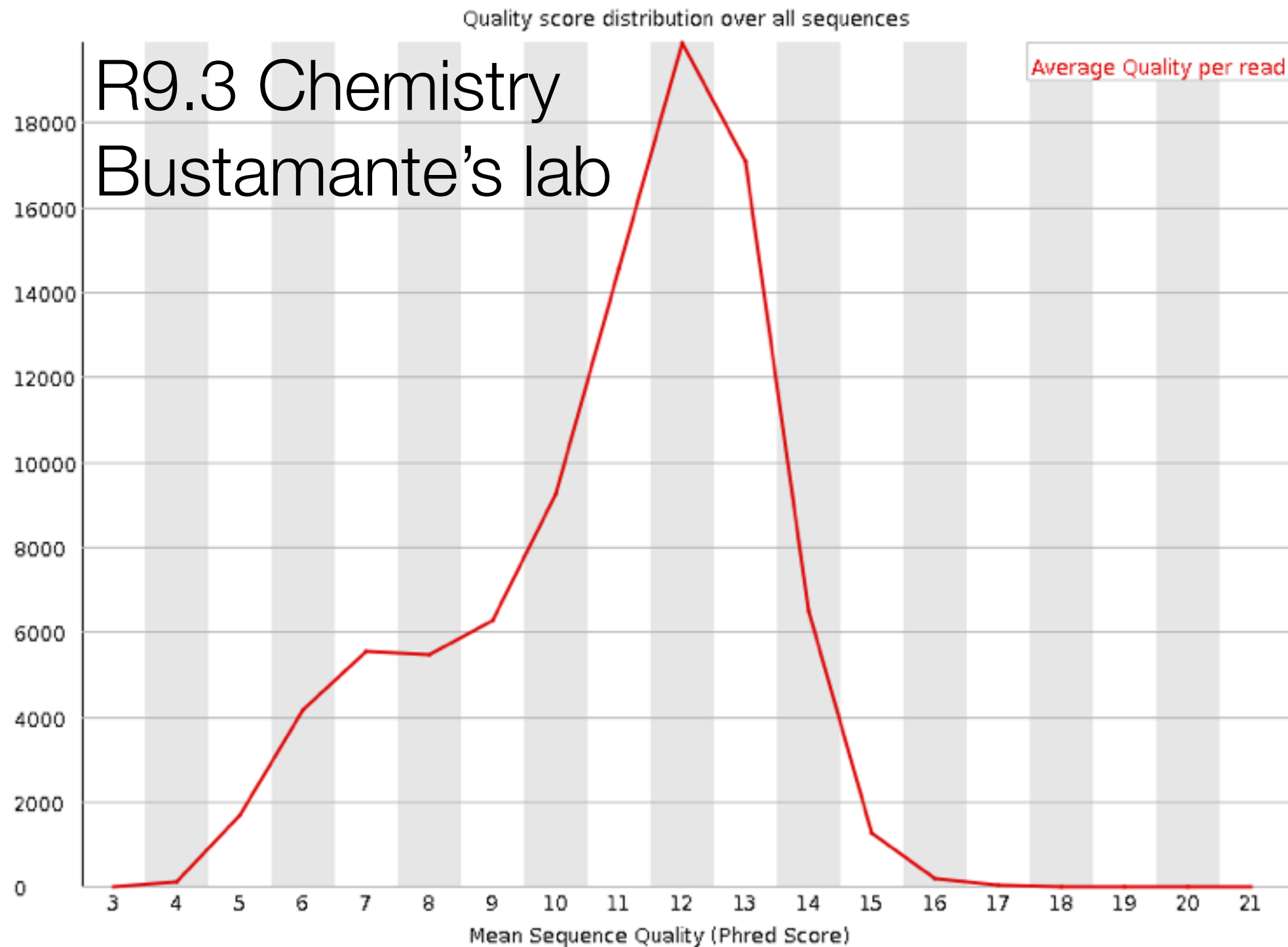
2D reads – typical profile of lambda DNA
mechanically sheared to average of 8 kb

1D reads – typical profile of lambda DNA
with no shearing

3

# Error rate is around 10%: Quality score freq. dist.



R9.3 Chemistry
Bustamante's lab

# Related research questions for long-read data

- How can we de-noise long-read data

  - Can we learn error profile?

  - Can we fix homopolymers?

    - Improve event calling step in base calling


- Variant calling from collection of mapped reads

  - How can we improve accuracy?

# [ Resources ]
# Long read: analysis tools & data sets

- Tools

  - <u>poretools</u>: **sequence extraction** & statistics (partial support for R9.4)

  - <u>porekit</u>:     statistics & plots (partial support for R9.4 chem.)

  - <u>minoTour</u>:   real-time data analysis (complicated installation)

  - <u>nanopolish</u>: compute consensus by HMM (de-noising)

  - <u>AlignQC</u>:     quality assessment of alignment (error model)

  - <u>DeepNano</u>:  base caller with Deep RNN (overfitting to a specific dataset?)

- Data sets (nanopore MinION & PacBio)

  - <u>NA12878 on MinION</u> by <u>Oxford Nanopore</u>

  - <u>NA12878 on MinION by Wellcome Trust Centre for Human Genetics</u>

  - <u>NA12878 on PacBio</u> by <u>Genome in a Bottle Consortium</u> (Mt.Sinai)

6

# My project
# Long read sequencer & Haplotype reference panel
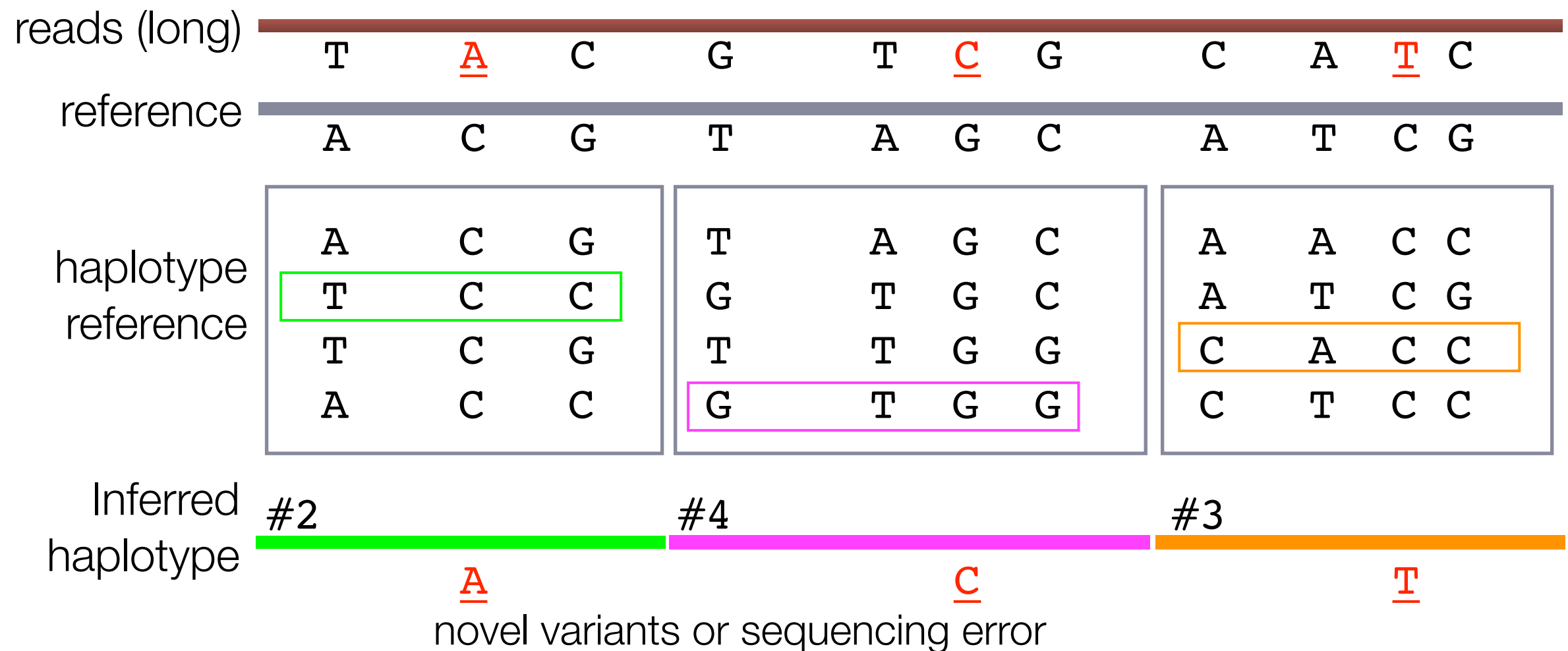
- Long read sequencers (Oxford Nanopore, PacBio)

  - Read length >= 8 kb

  - <u>Review paper of Nanopore</u>

- Haplotype reference

  - UK BioBank (152,729 imputed haplotypes)

  - Haplotype reference consortium (not available yet?)

- Compressed data representation of haplotypes

  - PLINK2

# [ Research Question ]
# Compressed representation of personal genomes

- How can we compress personal human genome ?



novel variants or sequencing error

- Indices for haplotypes + novel variants

8

# Compressed representation of personal genomes

- Validating the method with NA12878

- How can we discriminate variants from sequence errors?

- Is it possible to infer haplotypes on the fly?

- Is it scalable to a large population (say 100M)?

- Examples of applications:

  - Haplotyping for forensic science

  - Haplotyping on HLA region

  - Clinical applications