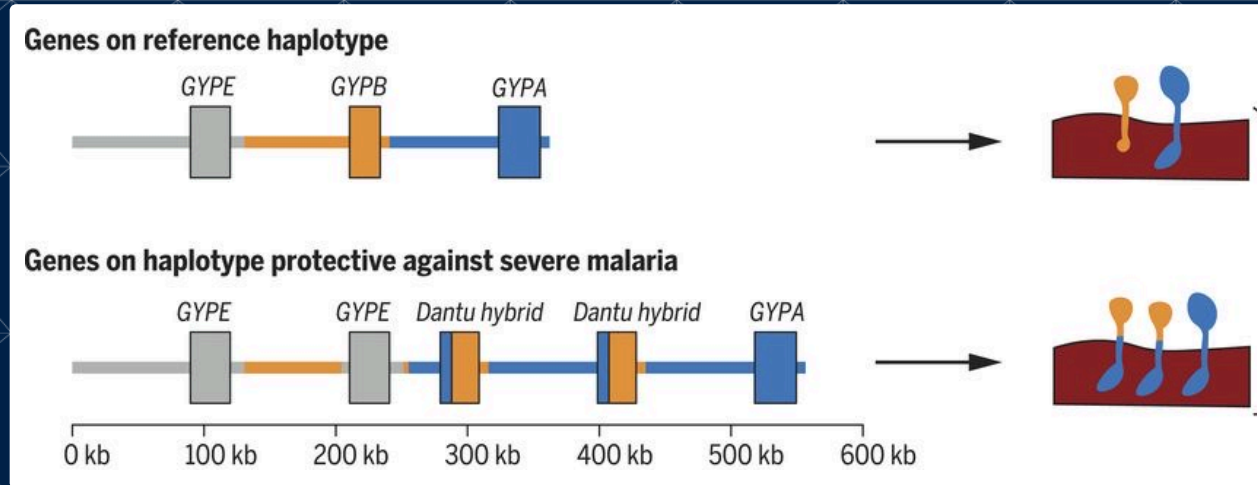# Long read accuracy and genome assembly

Gavin Band

LR CASe Detectives meeting
Septembr 7th 2023
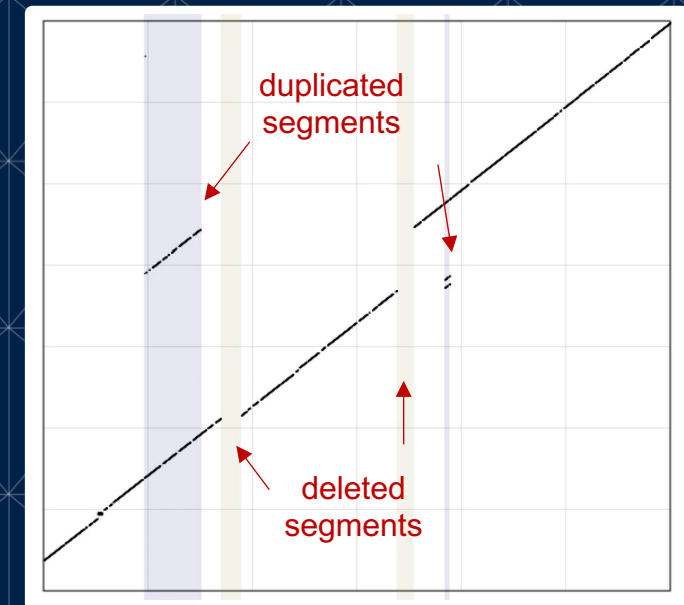
# Motivation: structural variation in hosts and pathogens



Genes on reference haplotype
GYPE  GYPB  GYPA

Genes on haplotype protective against severe malaria
GYPE  GYPE  Dantu hybrid  Dantu hybrid  GYPA

0 kb  100 kb  200 kb  300 kb  400 kb  500 kb  600 kb

"*Resistance to malaria through structural variation of red blood cell invasion receptors*"
2017

"*Malaria protection due to sickle haemoglobin depends on parasite genotype*"
2021



duplicated segments

deleted segments

Another *P.falciparum* genome

*P.falciparum* reference genome

# Many questions

- What is the structure of the variant?

- What is their functional effect?

- What is their phenotypic impact?

- How are they evolving?

- What other variants segregate?

-  How can we genotype them?

# Talk outline

1. How accurate are recent long-read platforms?

2. Two genome assembly applications
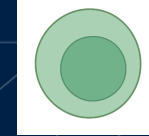
# The HV31 omniome project: data

## Genomic data:

- Illumina and MGI short-read data, to ~200x

- PacBio 'continuous long reads' (Sequel II), to ~35x
- PacBio 'HiFi' reads (Sequel II and IIe) to ~24x
- **New!!** PacBio 'HiFi' reads (Revio), to ~57x

- Oxford Nanopore Technologies R9.4.1, to ~63x
- **New!!** ONT R10.4.1 data, to ~69x

- 10X linked-reads (to ~40x)
- MGI stLFR linked reads

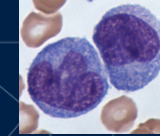- BioNano optical mapping, to ~150x coverage by fragments
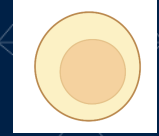
## Functional data:
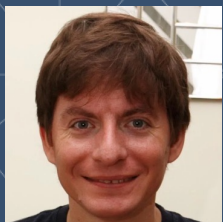


(B cells)    (T helper cells)    (Monocytes)    (Cytotoxic T)

- RNA-seq for gene expression
- ATAC-seq for chromatin accessibility
- CHiP-seq for histone modifications

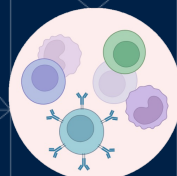- Methylation (from long read datasets)

All data is, or will be available through the EGA: **EGAS00001005046**

# Read length comparison
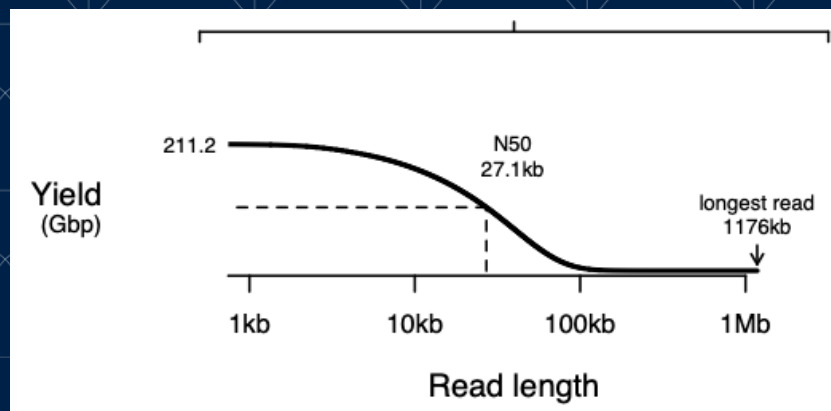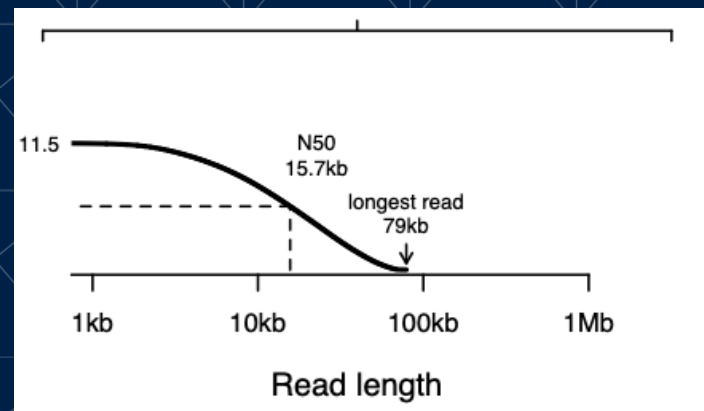
**Simplex** reads



Nanopore R10.4.1

In this expt, most nanopore reads
were 1-100kb long…

# Read length comparison

**Simplex** reads
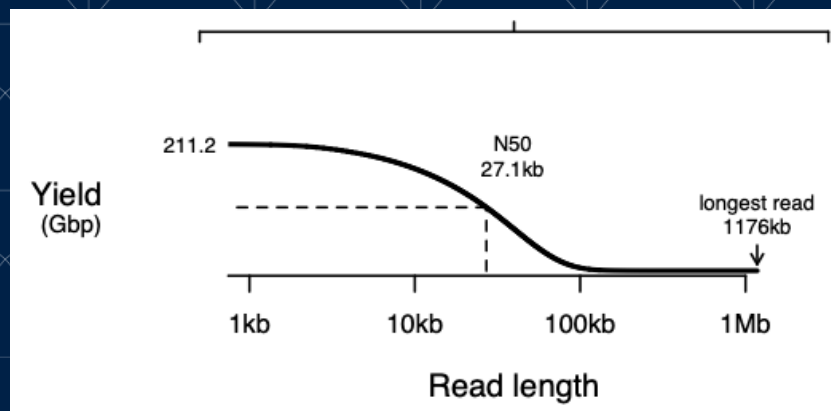
**Duplex** reads
about 5% of total reads

Nanopore R10.4.1
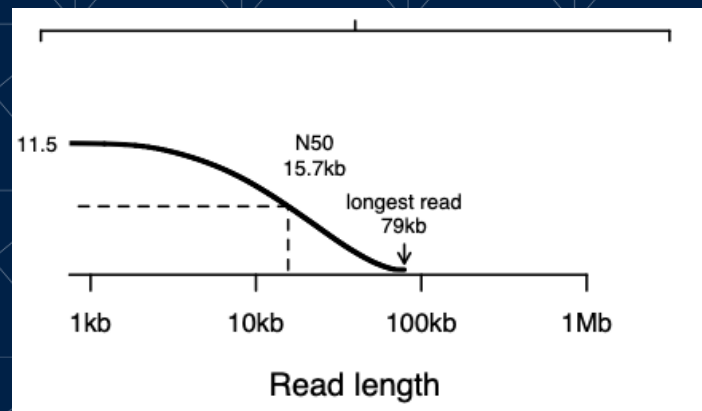
In this expt, most nanopore reads were 1-100kb long…

and duplex reads were slightly shorter
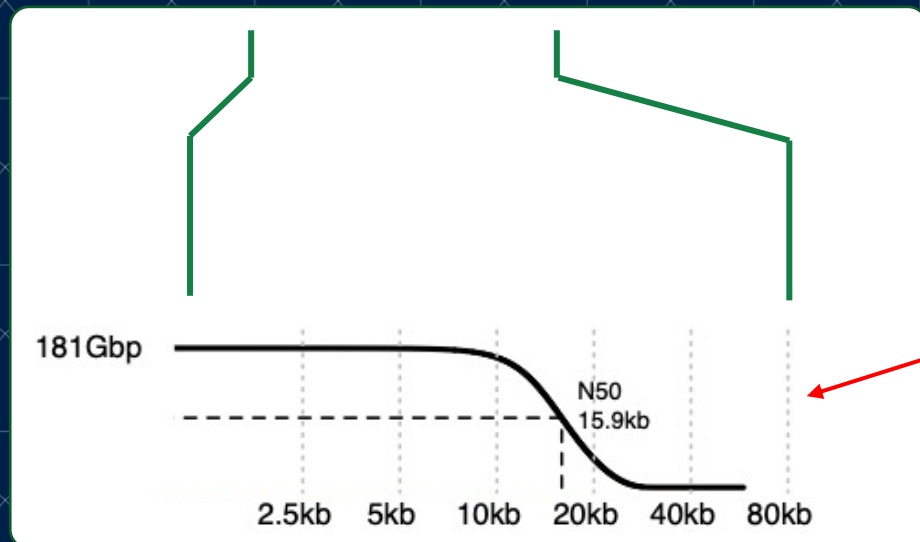
# Read length comparison

**Simplex** reads

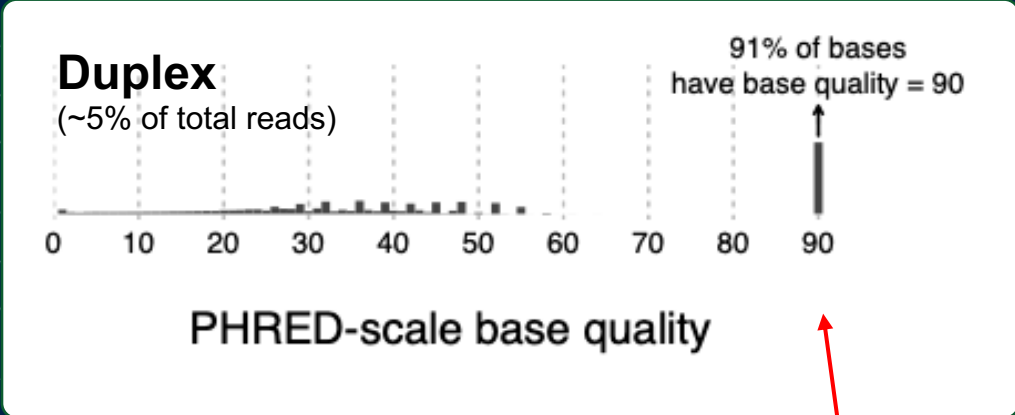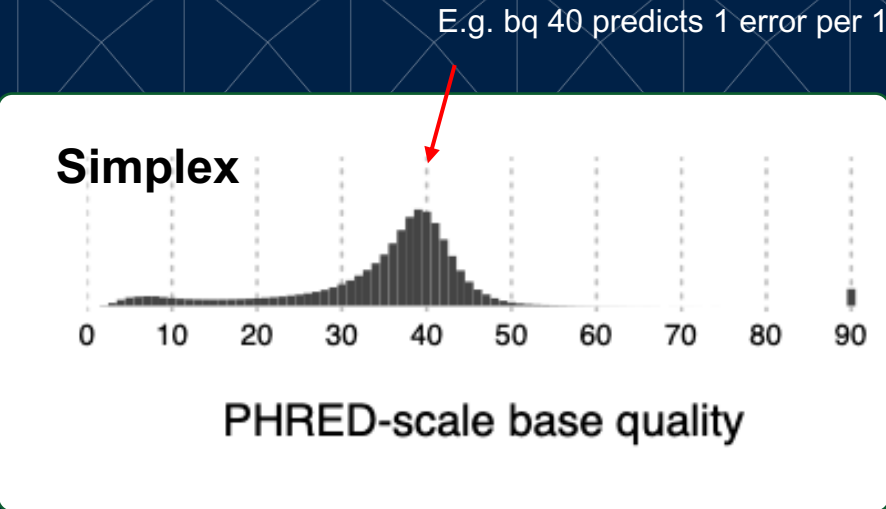**Duplex** reads
about 5% of total reads



Nanopore R10.4.1



Pacbio Revio

Pacbio reads are **shorter,** on average than nanopore reads.

About 10-20kb long

# Base quality comparison

E.g. bq 40 predicts 1 error per 10,000 bases

**Nanopore**
**R10.4.1**

**Simplex**

PHRED-scale base quality

**Duplex**
(~5% of total reads)

91% of bases have base quality = 90

PHRED-scale base quality

**Pacbio**
**Revio**

50% of bases have base quality = 40

PHRED-scale base quality

Pacbio base qualities are similar to nanopore simplex but compressed into discrete set of values to make the files smaller.
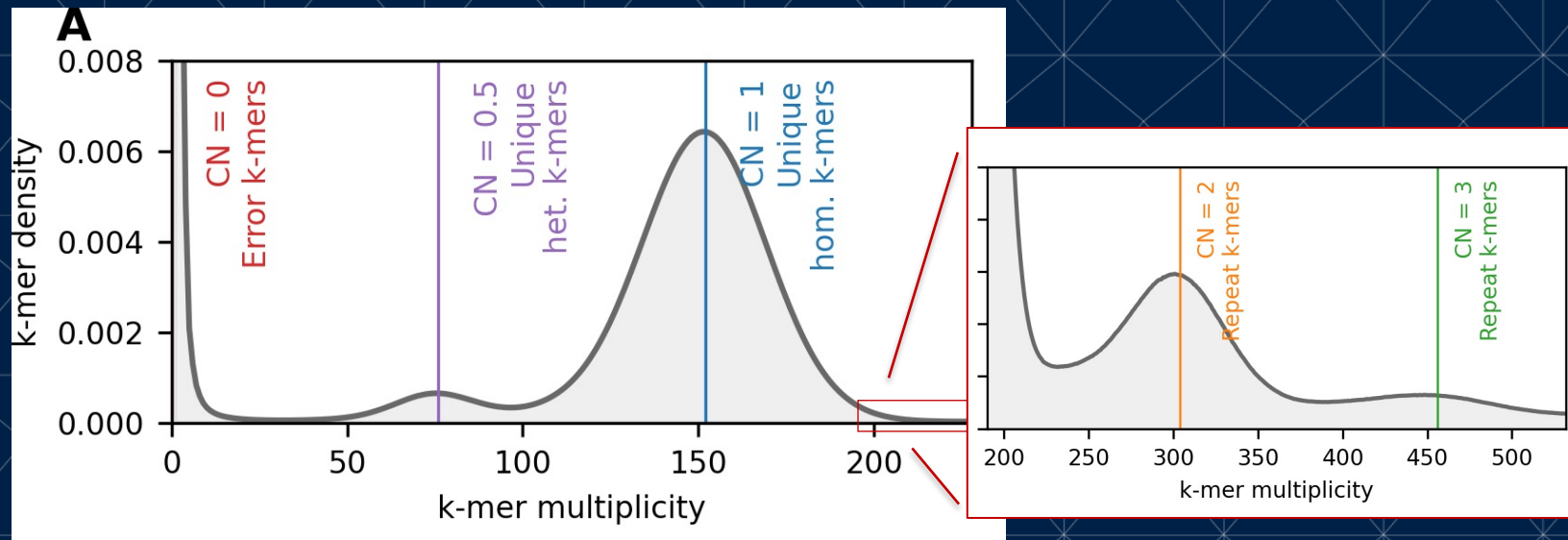
!

bq 90 predicts 1 error per billion bases

# Two ways to measure error rates

1. Measure *kmer accuracy* using a set of known true kmers

2. Measure base accuracy based on alignment to a reference
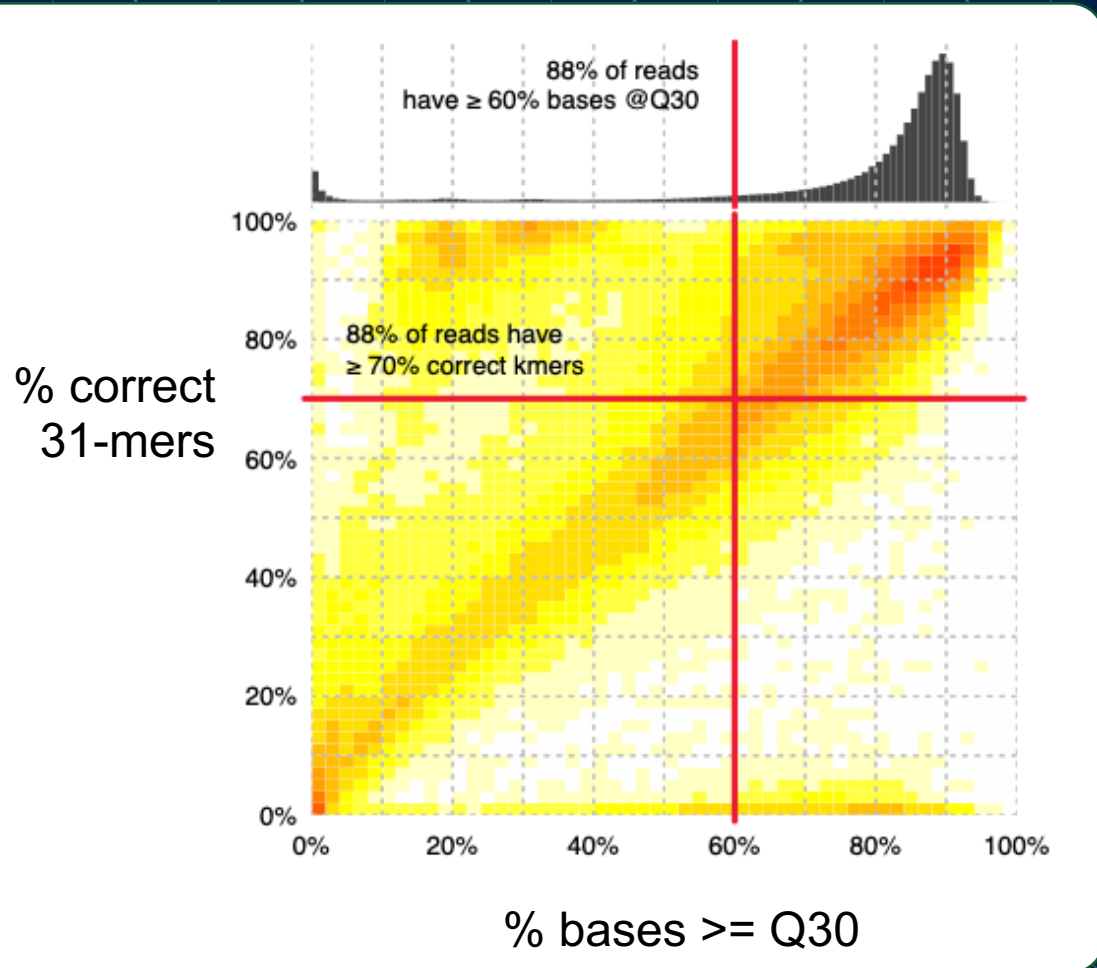
# Measuring kmer accuracy

**Method**: learn the set of true HV31 k-mers from short reads…



Histogram of k-mer multiplicity observed in Illumina, MGI, MGI CoolMPS, 10X and Sequel II data.

…and for each long read, count the number of true HV31 kmers (k=31)

* This histogram is mainly based on short reads, but also includes the 2019 Sequel II data.

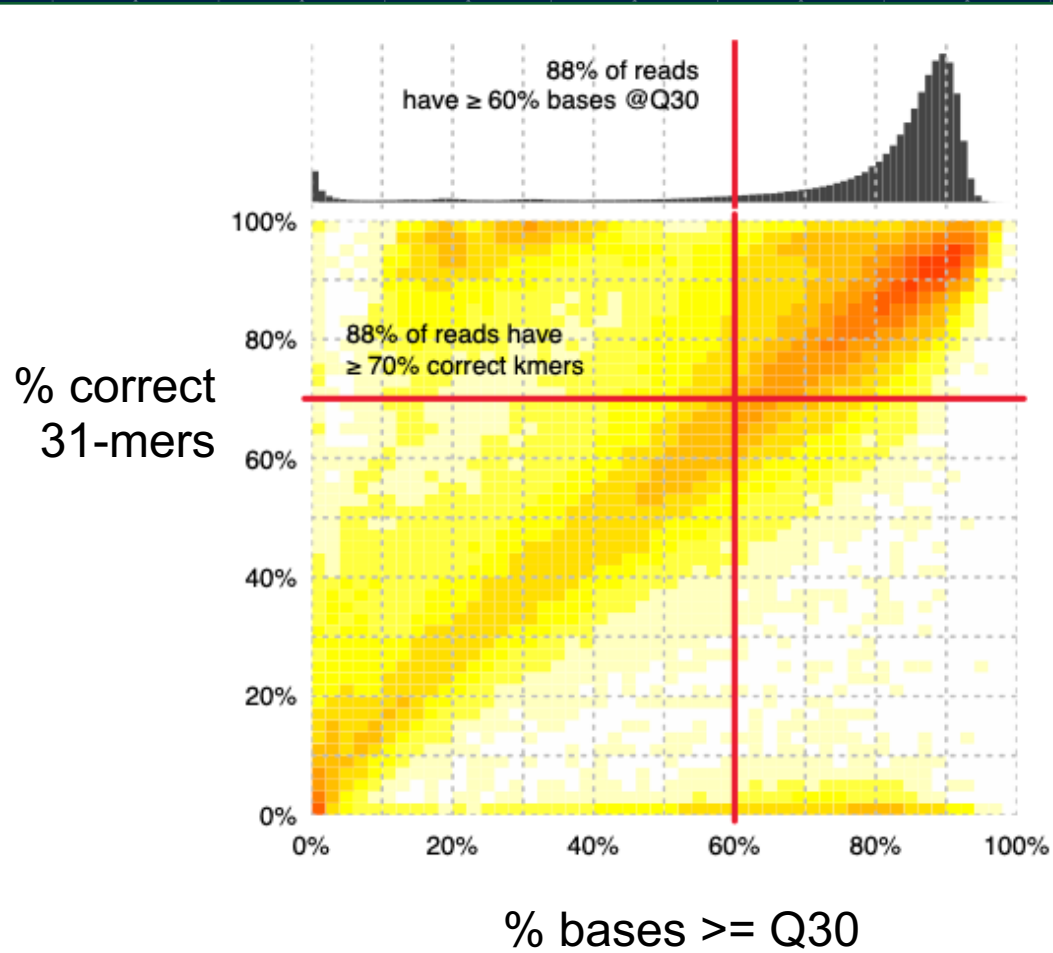# K-mer accuracy vs. predicted accuracy



**Nanopore R10.4.1 (simplex)**

Nanopore simplex has a roughly linear relationship between the quality predicted by base quality scores (x axis) and the observed quality (y axis)
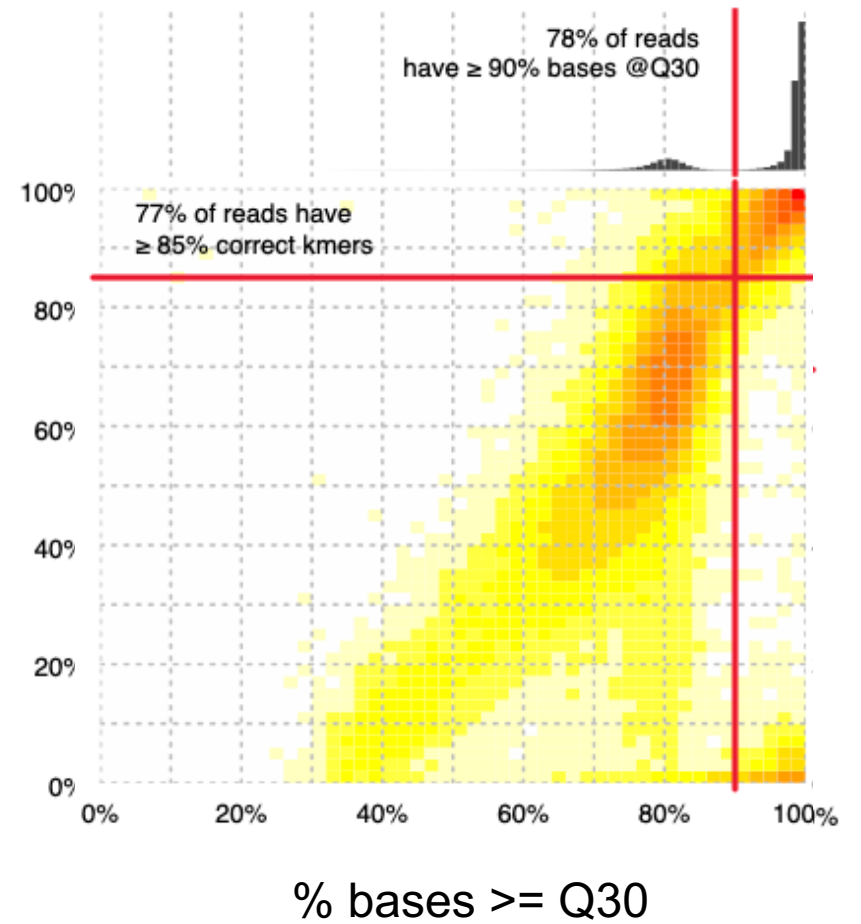
…as measured by accurate kmer rates.

Still about 20% of reads are poor quality.

# K-mer accuracy vs. predicted accuracy



**Nanopore R10.4.1 (simplex)**

**Nanopore R10.4.1 (duplex)**

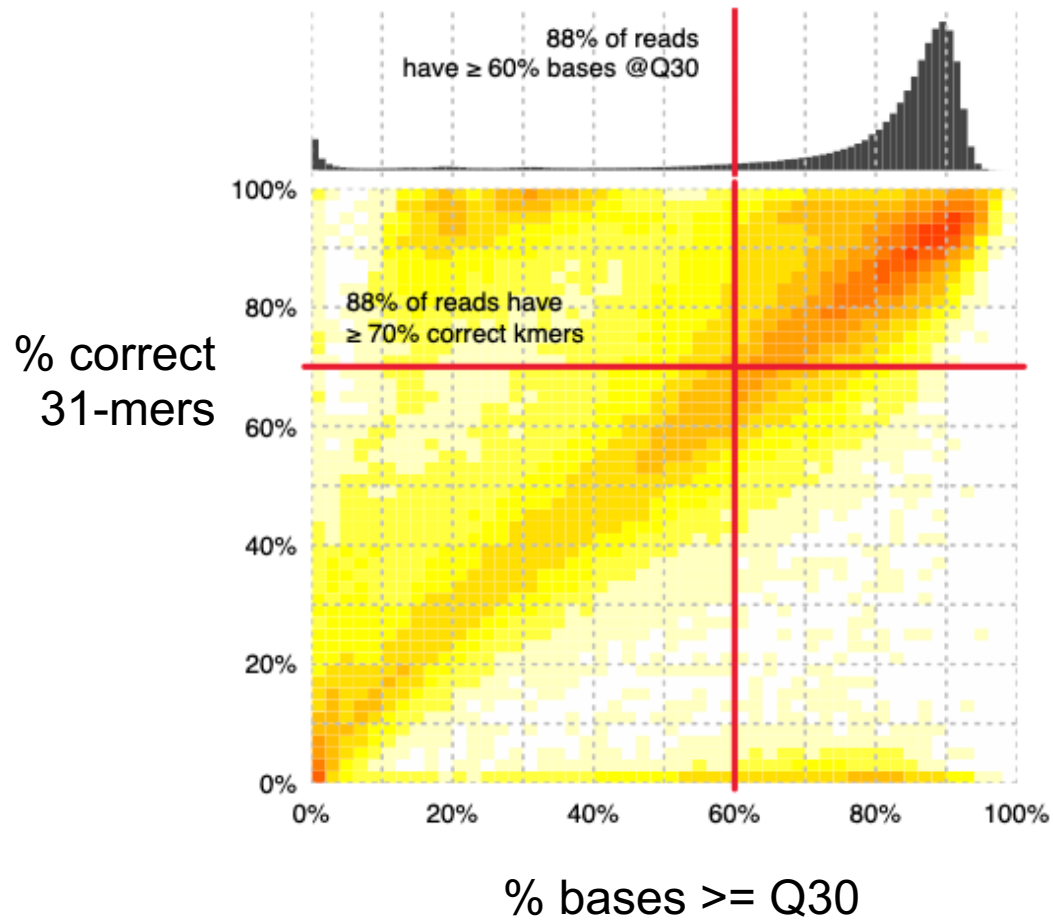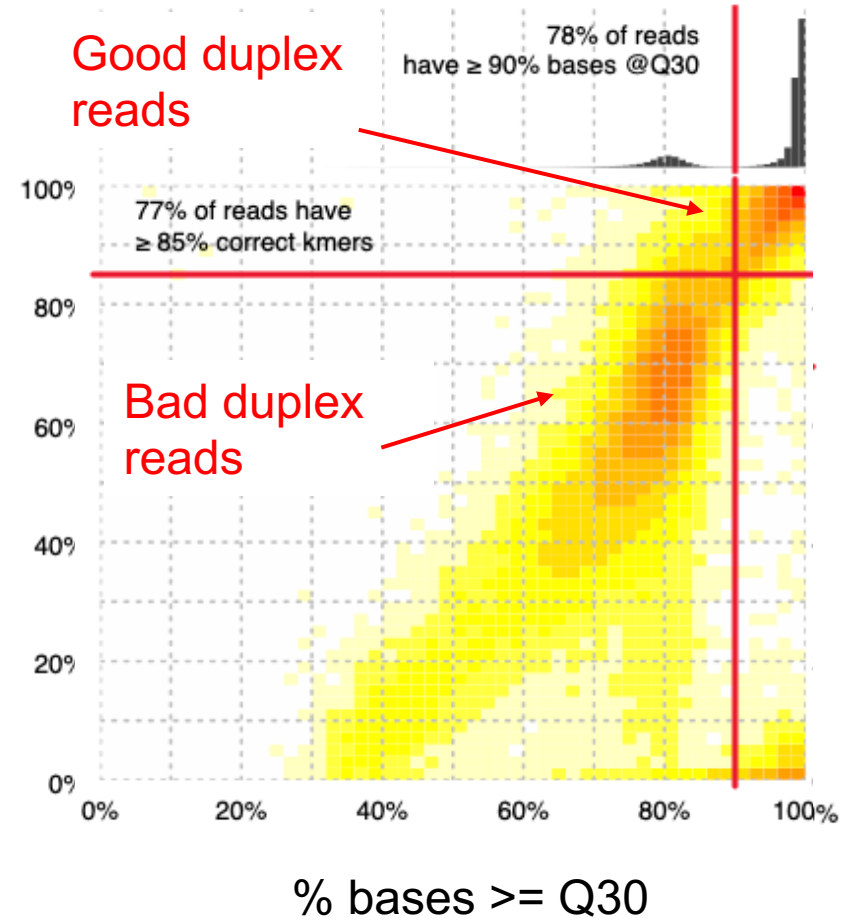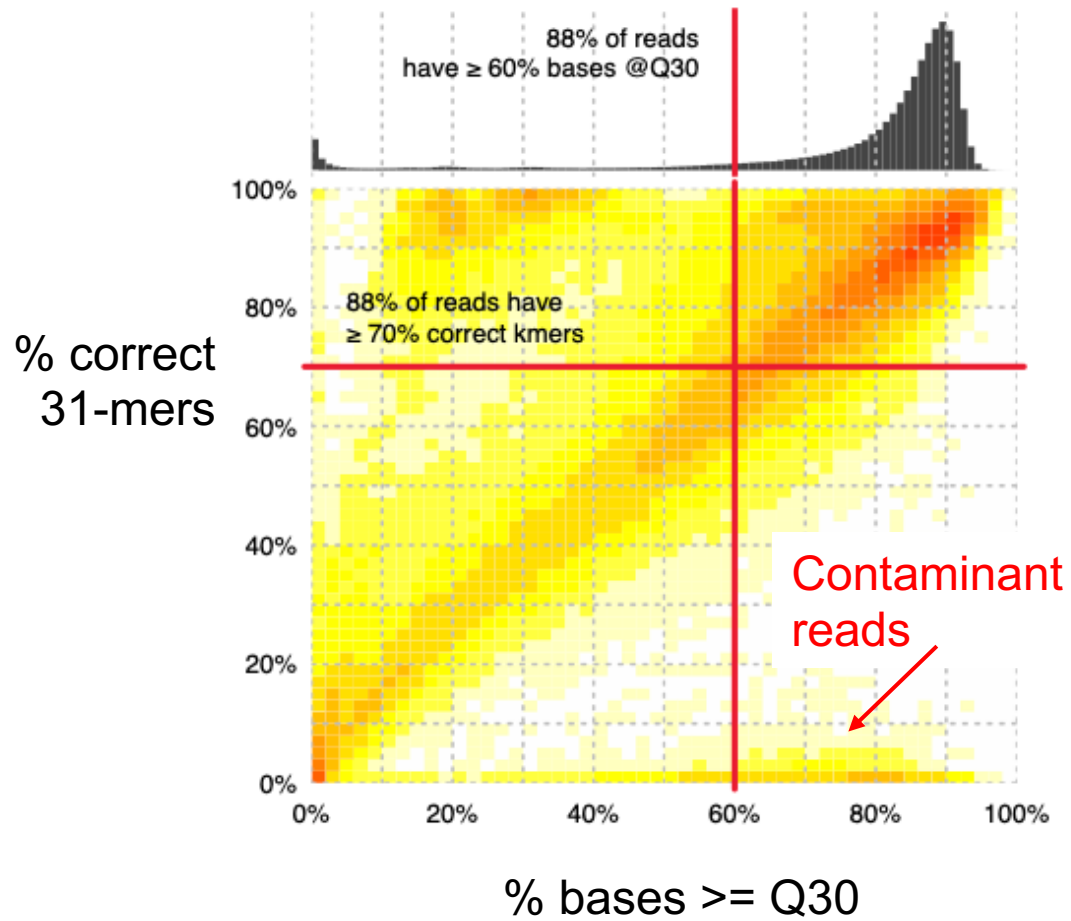# K-mer accuracy vs. predicted accuracy
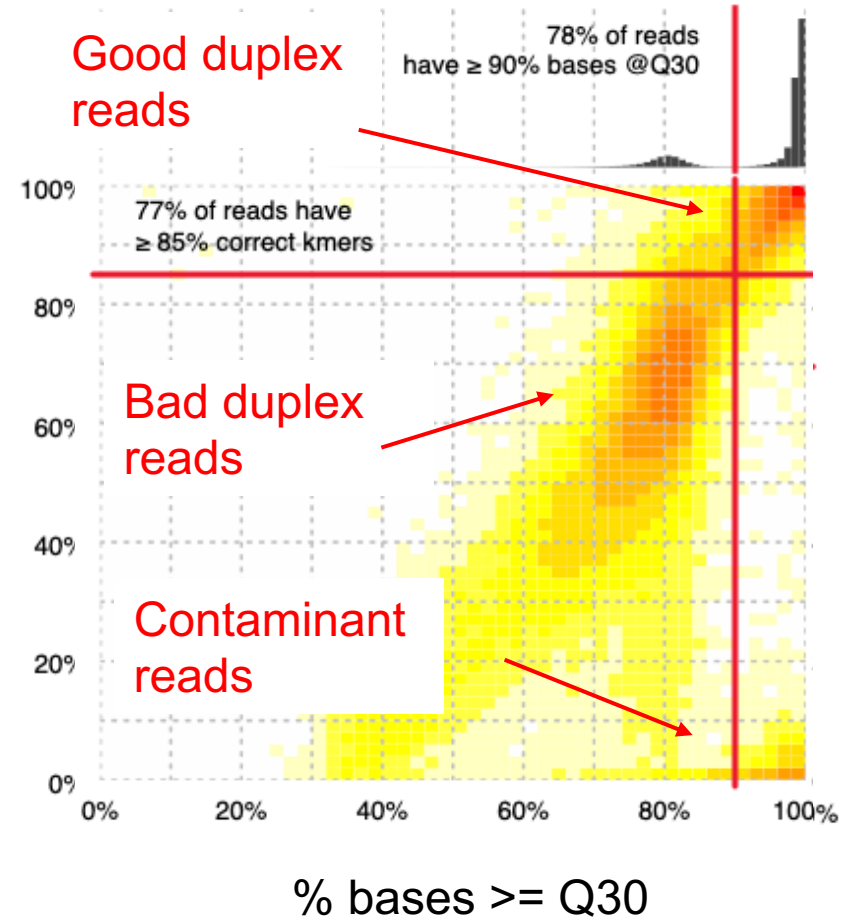


Nanopore R10.4.1 (simplex)



Nanopore R10.4.1 (duplex)

# K-mer accuracy vs. predicted accuracy



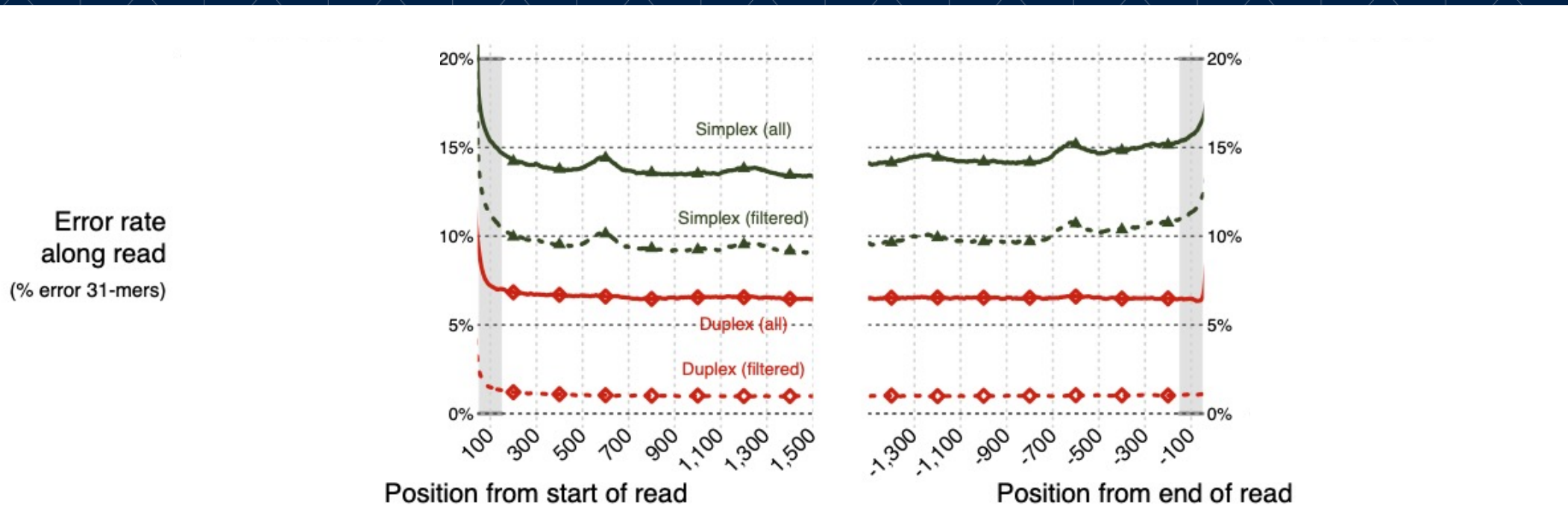Nanopore R10.4.1 (simplex)



Nanopore R10.4.1 (duplex)

# Accuracy along the read

Error rate along read

(% error 31-mers)

Simplex (all)

Simplex (filtered)

Duplex (all)

Duplex (filtered)

Position from start of read

Position from end of read
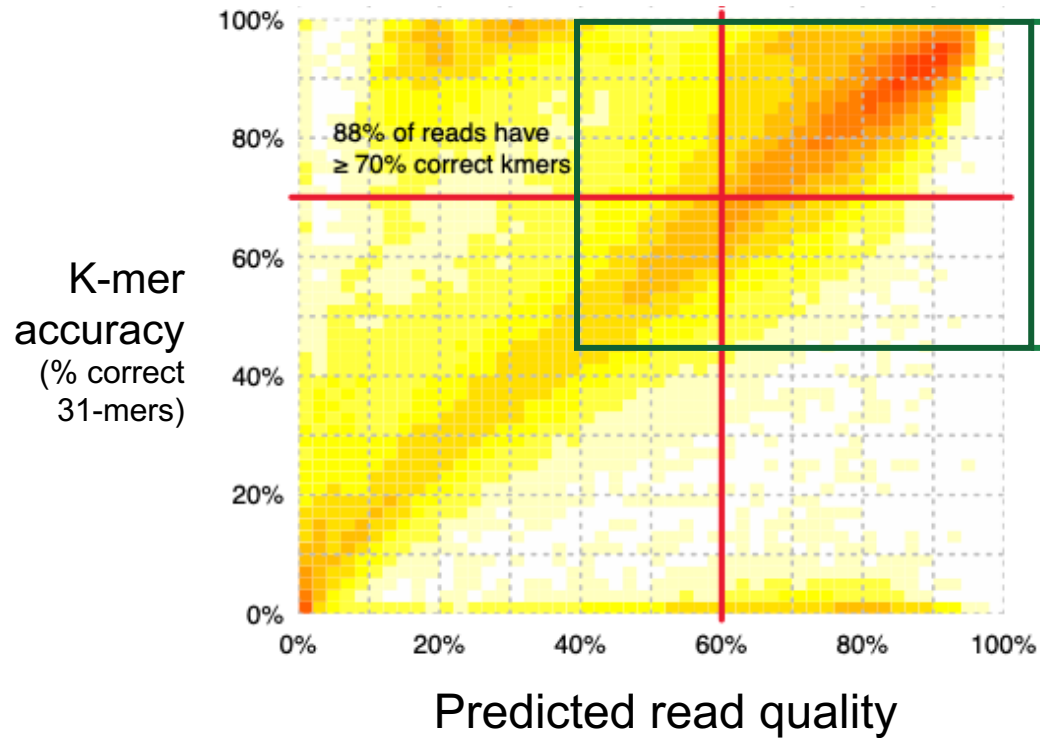
# Accuracy along the read



R10.4.1 is better than R9.4.1, especially after filtering.
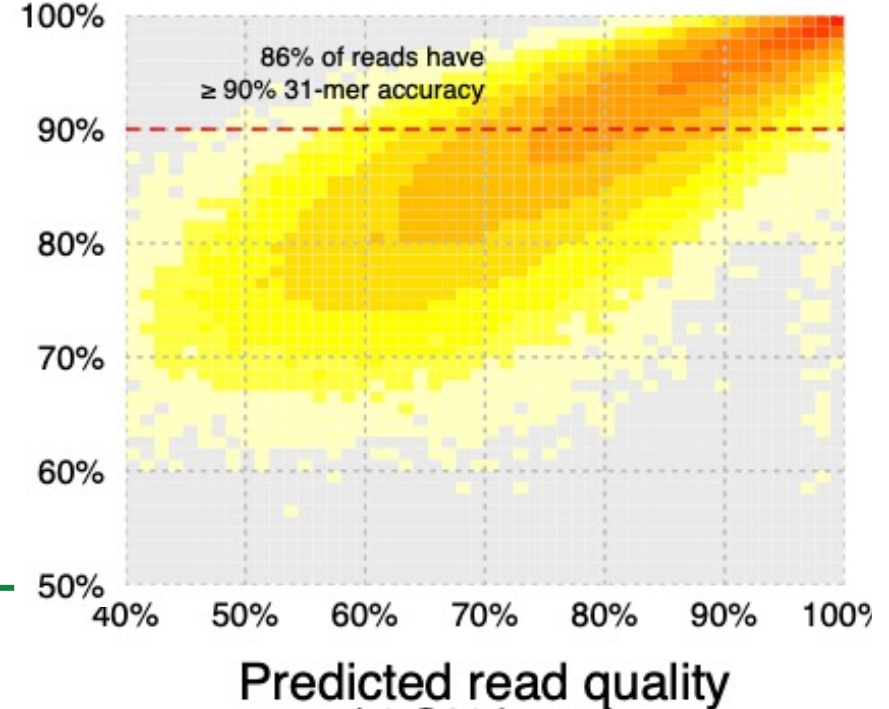Filtered duplex data has stupendously low error rates across most of the read.
Prominent read-end artifacts due to adapters (that might not be completely removable)

(Also, note the weird error bumps every 600bp…)

# K-mer accuracy vs. predicted accuracy



K-mer accuracy (% correct 31-mers)

88% of reads have ≥ 70% correct kmers

Predicted read quality

**Nanopore R10.4.1**
**simplex**

86% of reads have ≥ 90% 31-mer accuracy

Predicted read quality

**Pacbio Revio**

WELLCOME CENTRE *for* HUMAN GENETICS

Measuring predicted quality as: % bases >= Q30

# K-mer accuracy vs. predicted accuracy



K-mer accuracy (% correct 31-mers)

Predicted read quality

77% of reads have ≥ 85% correct kmers

86% of reads have ≥ 90% 31-mer accuracy

**Nanopore R10.4.1**
**duplex**

**Pacbio Revio**

Measuring predicted quality as: % bases >= Q30

# Accuracy along the read
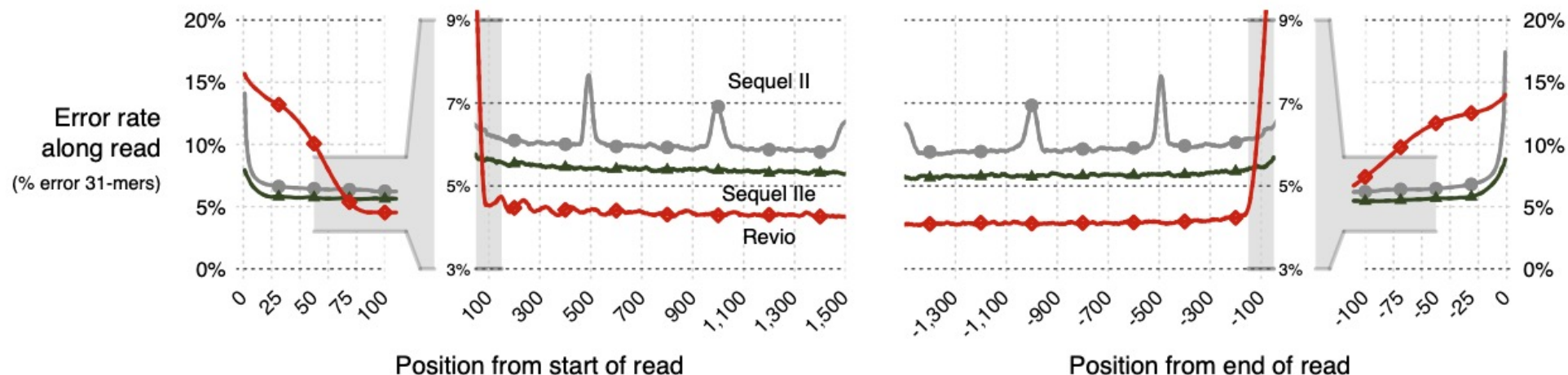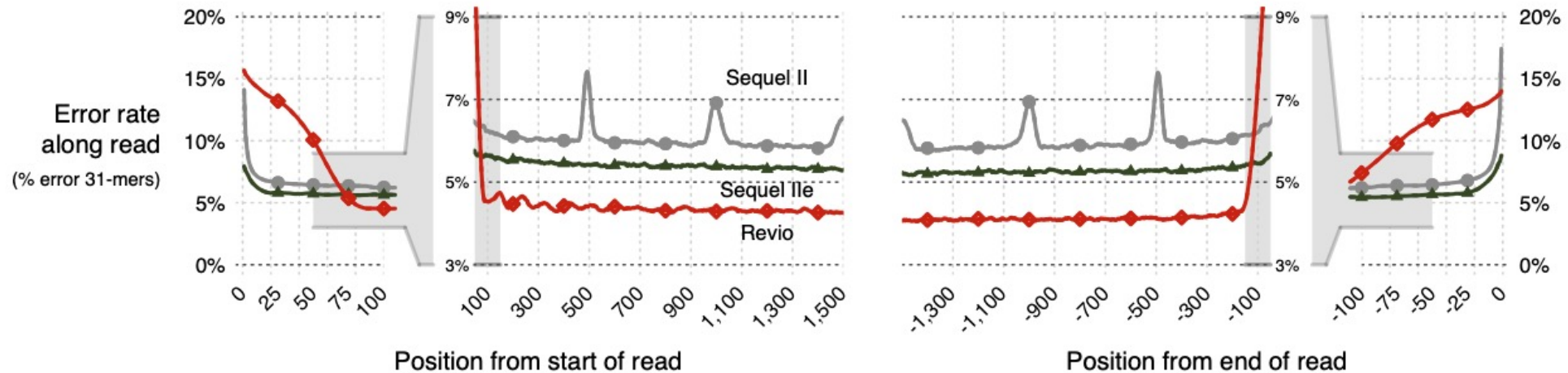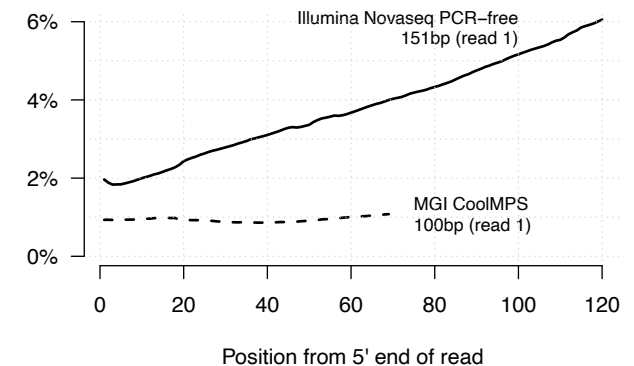


Our Revio data **also** shows elevated rates at the end of reads - !!
But improves upon Sequel IIe across most of the read length

(Meanwhile our older Sequel II data has weird, unexplained 'bumps' every 500bp.)

# Accuracy along the read



Error rate along read (% error 31-mers)

Sequel II
Sequel IIe
Revio

Position from start of read

Position from end of read

Error rates comparable to some Illumina data
Though some short-read datasets are better



Illumina Novaseq PCR–free
151bp (read 1)

MGI CoolMPS
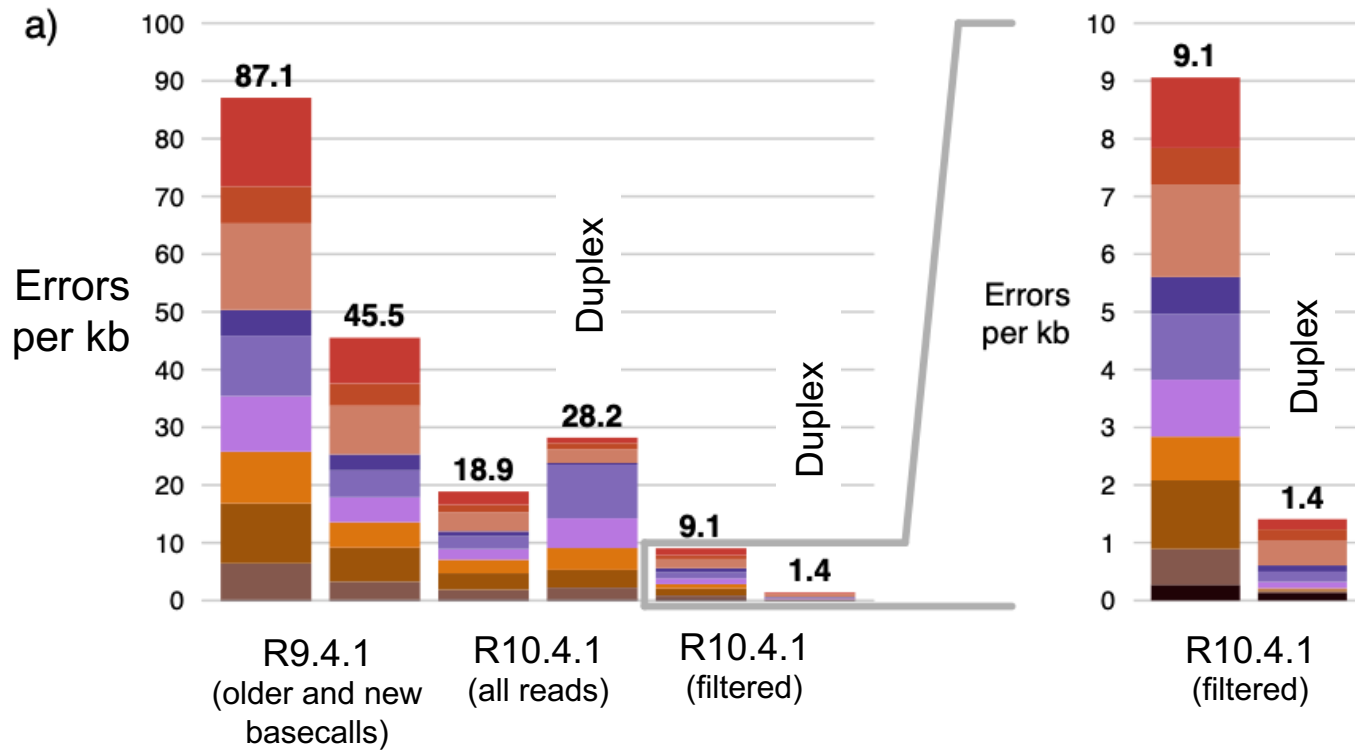100bp (read 1)

Position from 5' end of read

# Summary

- Nanopore R10.4.1 data improves over R9.4.1 data.

- Nanopore still noisy and has a few artifacts

- Pacbio Revio also improves over Sequel IIe across most of the read

- Nanopore duplex reads are somewhat comparable to Pacbio reads – maybe better after filtering, but are only 4-5% of data

- Both platforms have annoying-looking read-end effects.

Alternate approach: **align** to a reference sequence, **mask out** true variation and repetitive sequence

We use T2T assembly, mask out SNPs, INDELS, and SVs from HV31 data, and satellite arrays, segdups, repeat-masked elts.

# Nanopore

Nanopore

Pacbio

a)

Errors per kb

87.1

45.5

18.9

28.2

9.1  Duplex

1.4  Duplex

Errors per kb

9.1  Duplex

1.4

R9.4.1
(older and new basecalls)

R10.4.1
(all reads)

R10.4.1
(filtered)

R10.4.1
(filtered)

3.1

2.7

2.3
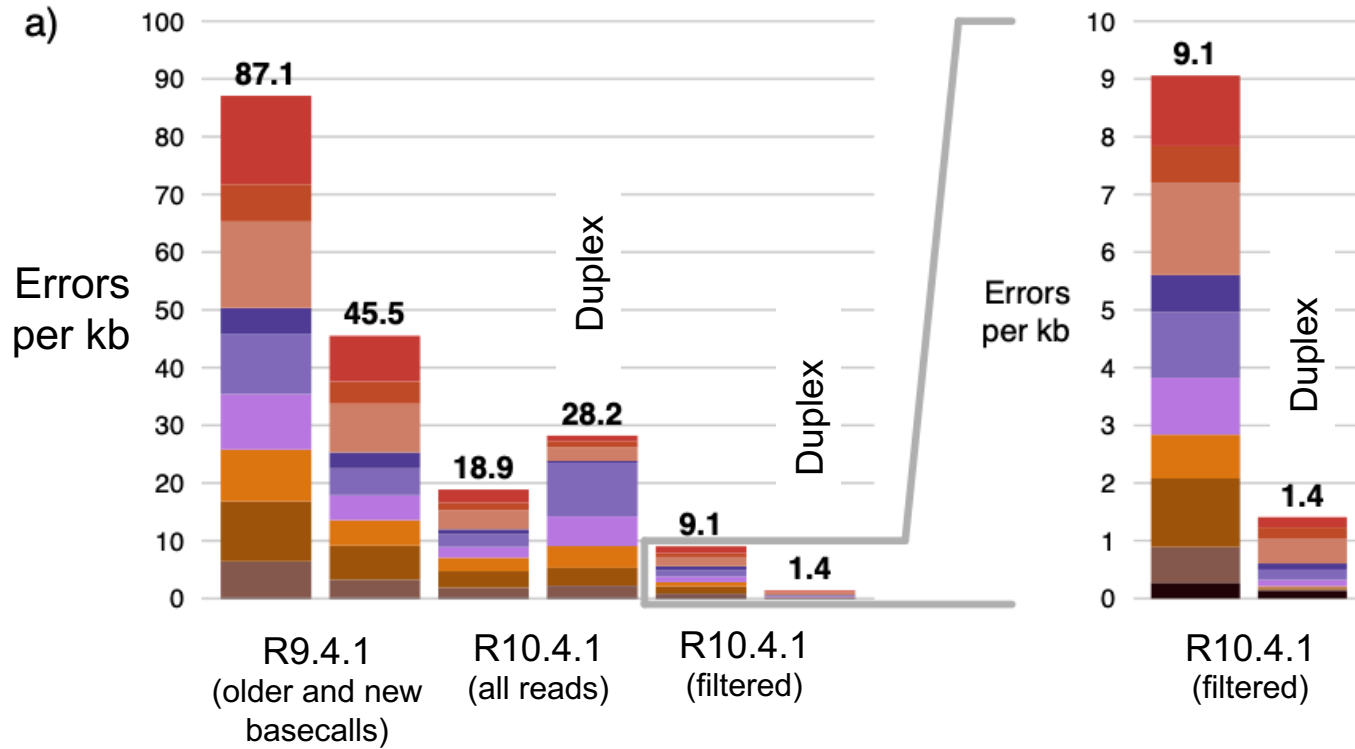
Sequel II

Sequel IIe

Revio

- Homopolymer contraction
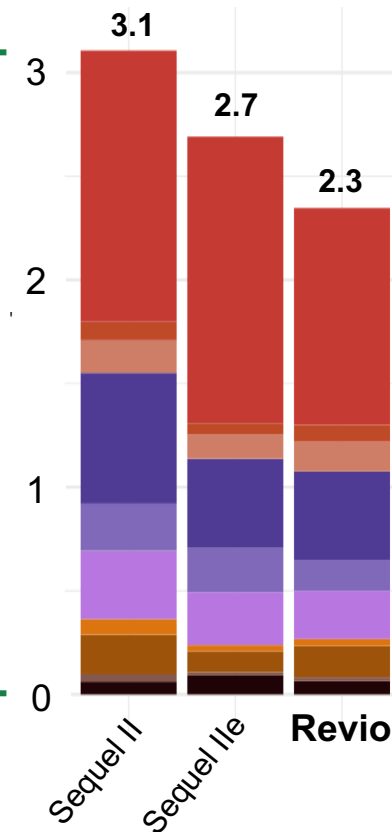- Homopolymer-creating deletion
- Other deletion
- Homopolymer expansion
- Other insertion in homopolymer
- Other insertion
- Substitution in homopolymer
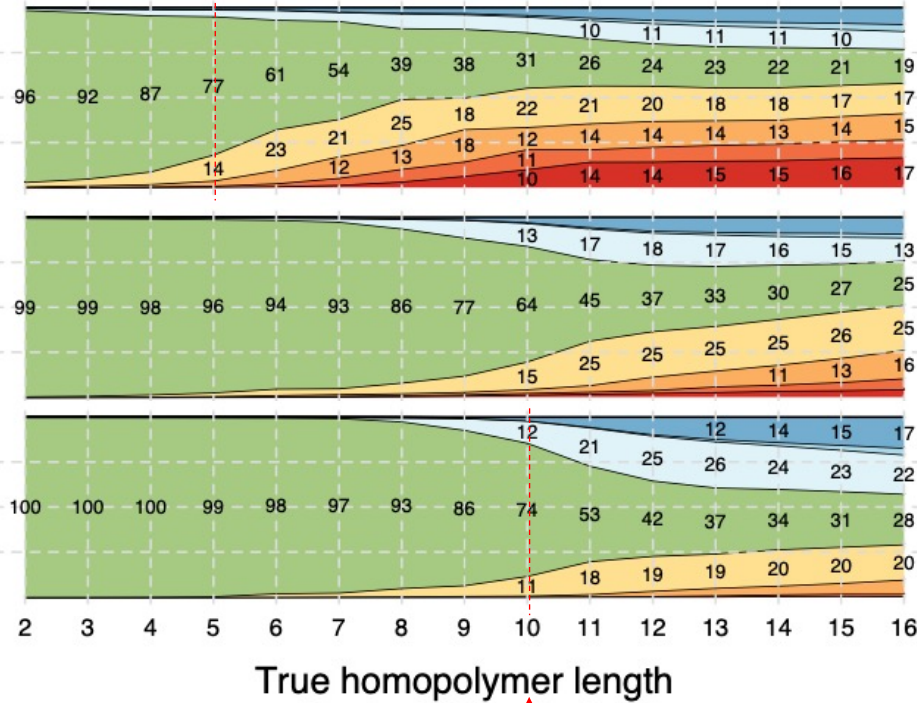- Homopolymer-creating substitution
- Other substitution

WELLCOME CENTRE for HUMAN GENETICS

25

R9 pore size

R10 pore size

Nanopore homopolymer length calling still drops off above pore size…

(but R10 pore size is larger)

Nanopore homopolymer length calling still drops off above pore size…

(but R10 pore size is larger)

Pacbio calls longer homopolymers better
still only ~60-70% accuracy for longest lengths

# Subtle substitution biases are also present



Nanopore tends to make transition-like errors
(A<->G and C <-> T).

CpG sites appear to have a particularly high substitution rates.
But the absolute rate is still low.

# Subtle substitution biases are also present
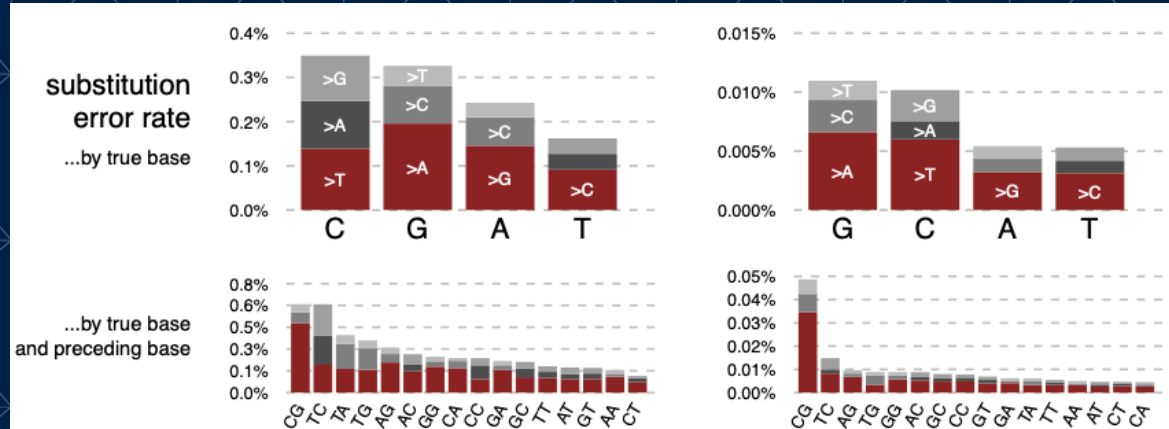
Simplex

Duplex



Nanopore makes substitutions of C and G bases, and tends to make transition-like errors (A<->G and C <-> T).

CpG sites appear to have a particularly high substitution rates.
But the absolute rate is still low.

Pacbio



Pacbio makes more substitutions at A and T bases and tends to miscall to A or T.

# Subtle substitution biases are also present



Moreover both platforms appear to have elevated substitution rates at CpG sites

# Summary

New revisions of ONT and Pacbio data are both fantastic.

Nanopore requires more downstream work to filter / process.

Duplex reads look very exciting, if low throughput can be overcome.

# Costs

For this experiment we 5 Promethion flowcells and 2 Revio SMRT cells were used.

For ONT, the list cost places the consumables cost at £2,700 - £4,050 flowcell cost, depending on order volumes, plus possibly £500 for library reagents. However you might only need 3 flowcells with current version (because it runs at a faster rate), so perhaps £2,120 - £2,930 in total

For Pacbio, it's a bit unclear to me but two flowcells might cost ~£2,000 with library prep on the order of £500 (I think - very ballpark.), so £2,500 in total.
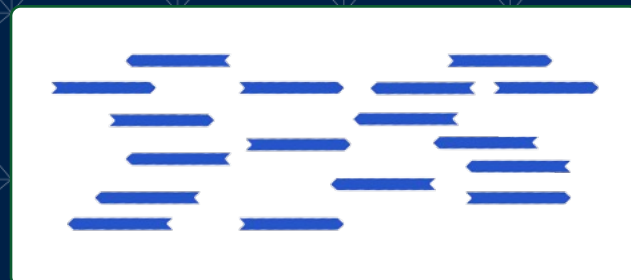
In other words - the costs look very similar to me.

Note these costs do **not** include equipment, service, personnel or additional reagent costs.

Genome assembly application 1

# A haplotype-resolved assembly with functional data



PacBio Sequel II/IIe
ONT R10.4.1

Jia-Yuan
Zhang

**RNA-seq**
(expression)

Verkko

Methylation

Align and
resolve phase

**ChIP-seq**
(For histone
modifications)

Multiple approaches
(BubbleGun, Linked reads,
kmer approach),
WhatsHap, HapCut2

Haplotype 1
Haplotype 2

**ATAC-seq**
(detects open
chromatin)

## Phased 'omniome"
reflecting immune cell types

Jia-Yuan Zhang

Example: a segmental duplication at TCAF1/2 locus
Not fully resolved in the Verkko assembly graph.

Use an empirical model of the k-mer distribution to probabilistically resolve the most-likely pair of haplotypes.

A ~50Mb 'phased' NG50
(50% of assembly bases are in phased contigs of 50Mb or greater)

Genome assembly application 2: resolving malaria structural variants involved in host-parasite interactions

# Three regions of the *Pf* genome are associated with sickle hamoglobin

Evidence for association
for *P. falciparum* variants
(averaged over human variants)

*P. falciparum* genetic variants

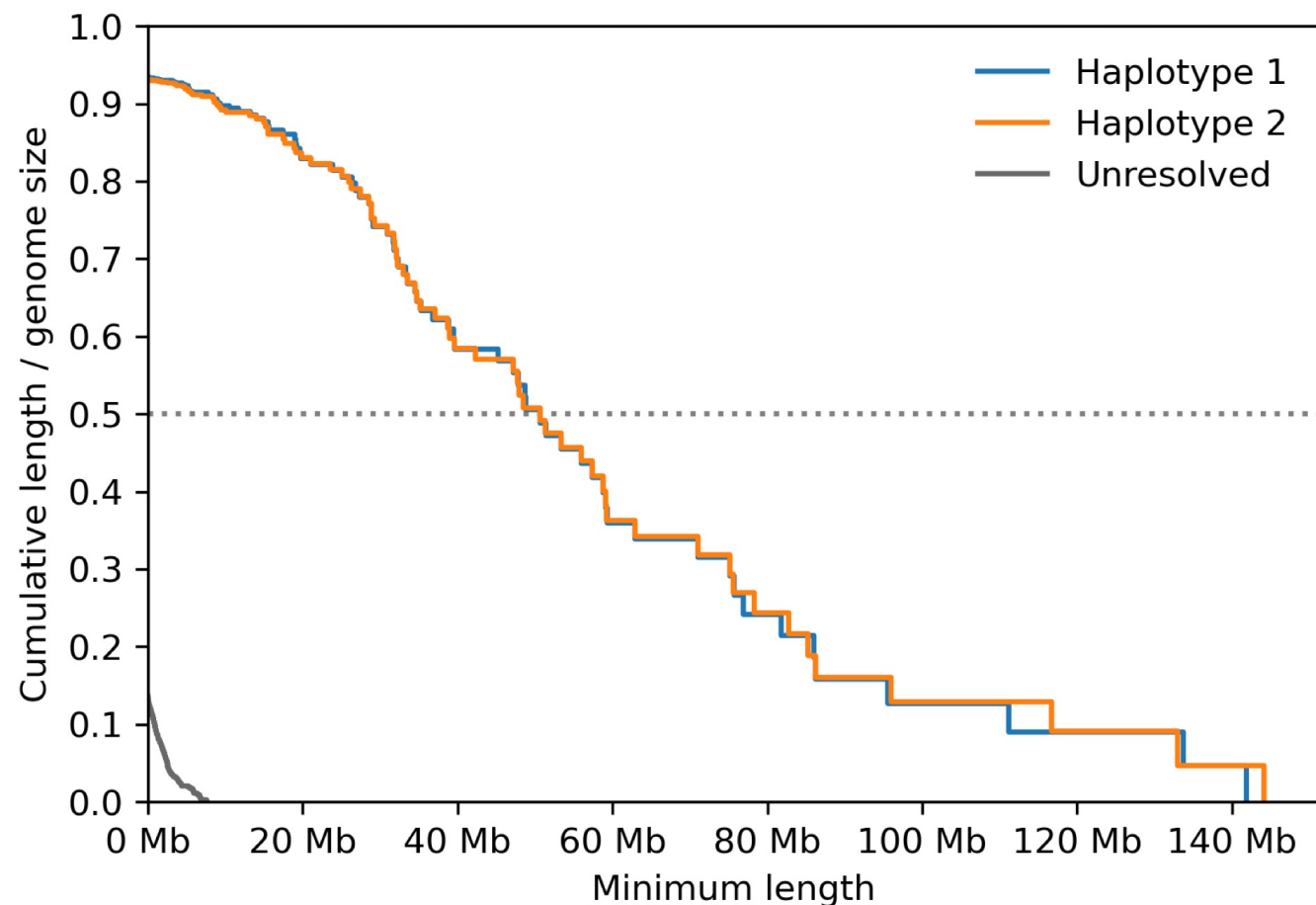HbS appears to give very strong protection against reference-like parasites, but maybe hardly any against + + + parasites

*Pfsa3*

*Pfsa2*

*Pfsa1*

Relative risk of HbS
by *Pf* genotype

*Pfsa* genotype
− = reference allele
+ = HbS–associated allele

$RR = 0.012$

$RR = 0.09$

$RR = 0.17$

$RR = 0.83$

Evidence for association
for *P.falciparum* variants
(averaged over human variants)

*Pfsa3*

Evidence for
association
for *P.falciparum*
variants
with HbS

chr11:1,058,035

The top SNPs are non-synonymous changes.
**However** they also appear to be linked to a surrounding structural variant, and are associated with increase transcription.

duplicated
segments

deleted
segments

*Reference parasite*

Attempt 1: Nanopore-based amplicon sequencing

Annie Forster

Jason Hendry

Mariateresa de Cesare

Anna Jeffresy

Analysis of short read data (MalariaGEN PF6) revealed there are multiple structural types segregating.

Kmer sharing in HbS-associated regions

04/11/2021, 10:44

**Annie Forster**

**Nanopore amplicon sequencing:**
Jason Hendry
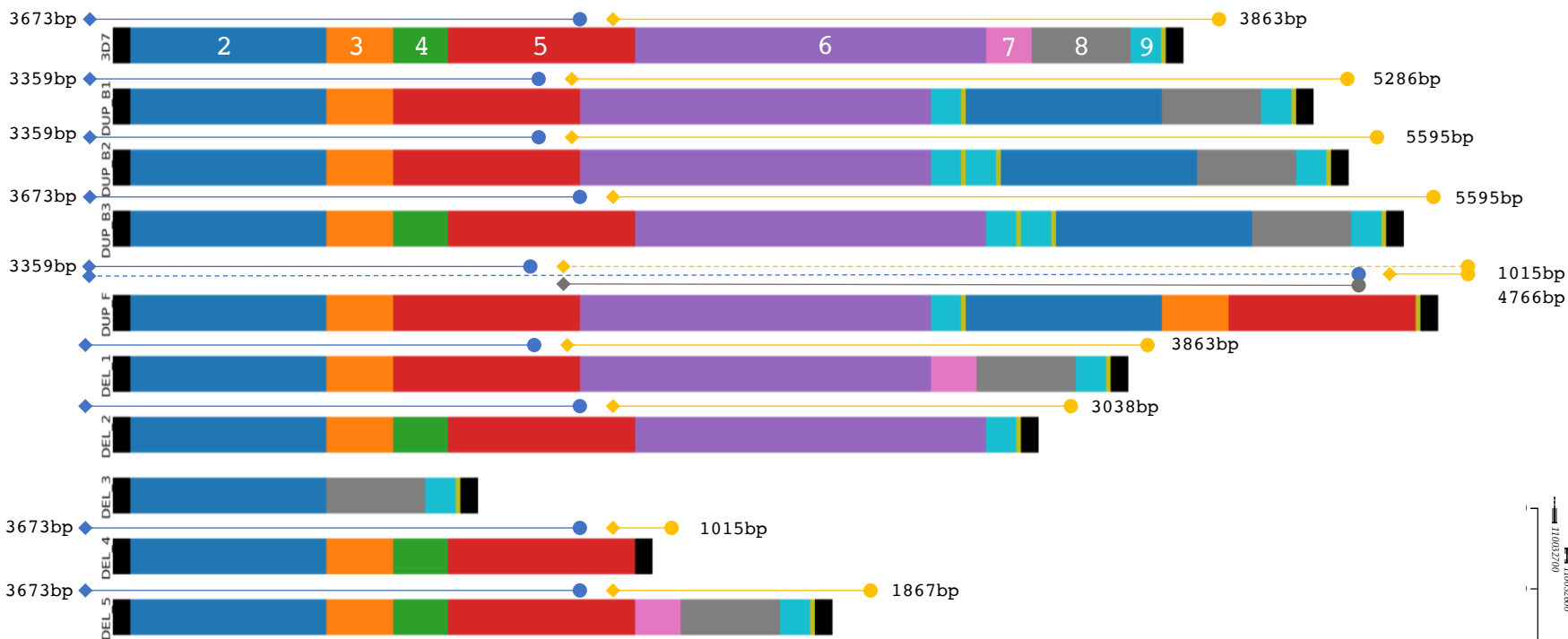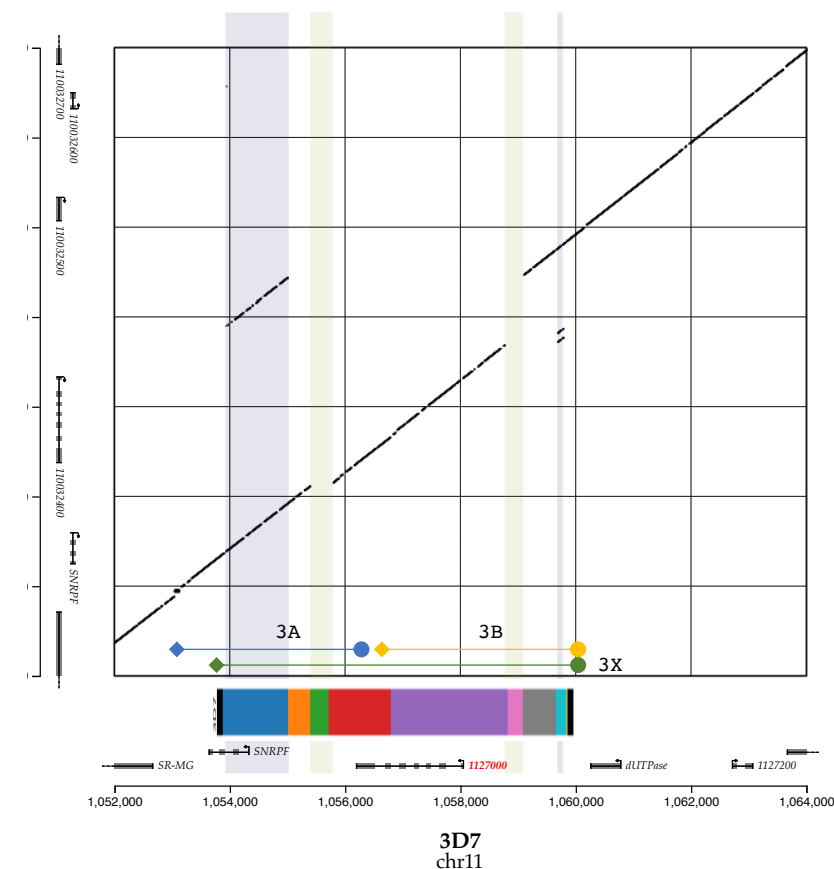Mariateresa de Cesare
Anna Jeffreys

3D7
chr11

FUP_multiplex01 aligned to a mock-up DUP_F reference. Looks like there are three fragments as predicted! It's a bit difficult to count length but roughly they seem to be...

1: 3350bp
2: 4651bp
3: 979bp?

Predicted lengths were:

Pfsa3A – 3,359bp
Hybrid – 4,766bp
Pfsa3B – 1,015bp

2nd attempt: Pacbio whole-genome sequencing

3D7
FUP-H
4 Kenyan parasites
2 Gambian parasites
1 parasite from single-cell sorting

Carried out by James Docker and Amy Trebes, Oxford Genomics Centre for a test of new fragmentation protocol.

Worked amazingly well

Alex Macharia, Patrick



Parasite from blood culture Kilfi

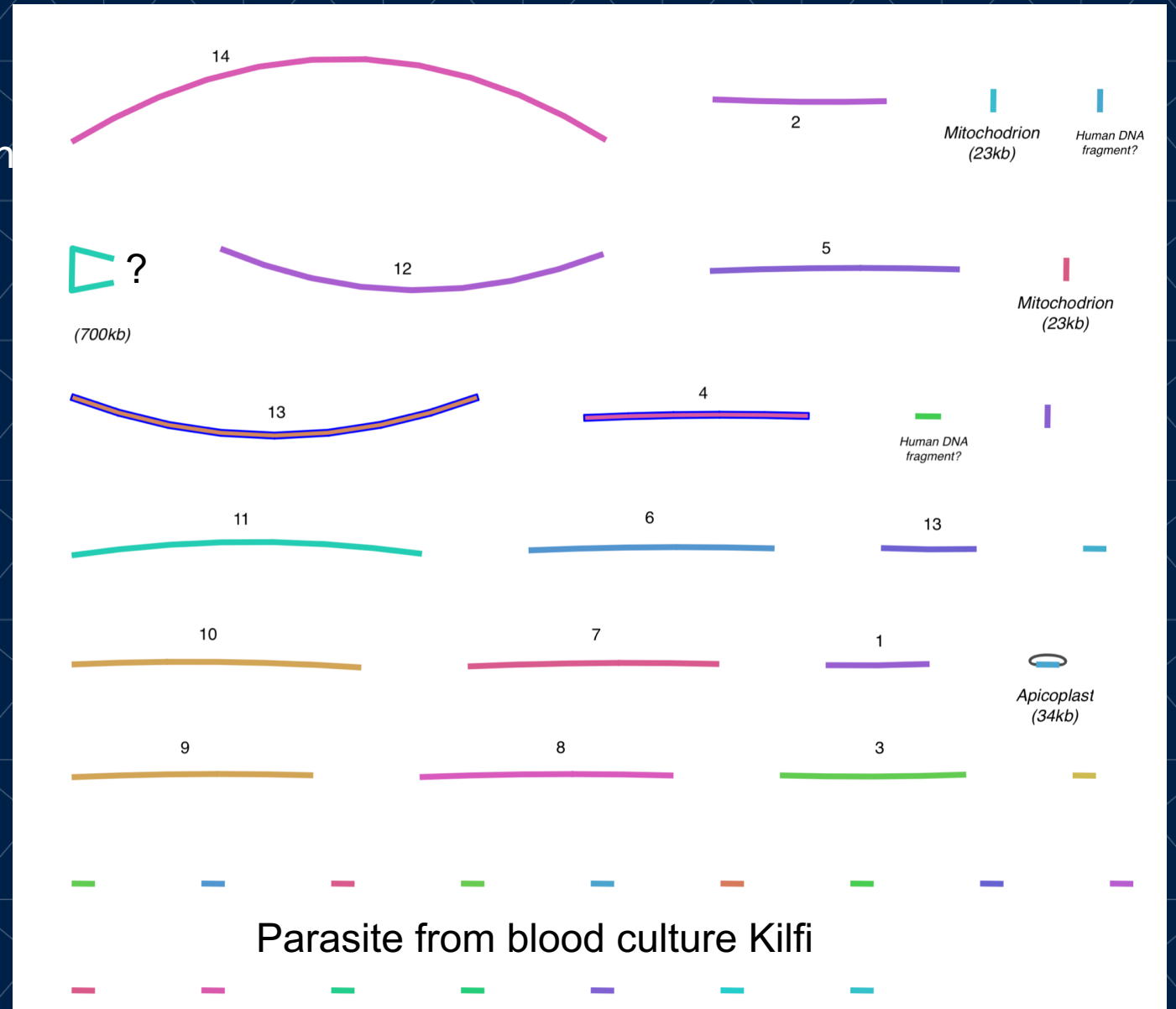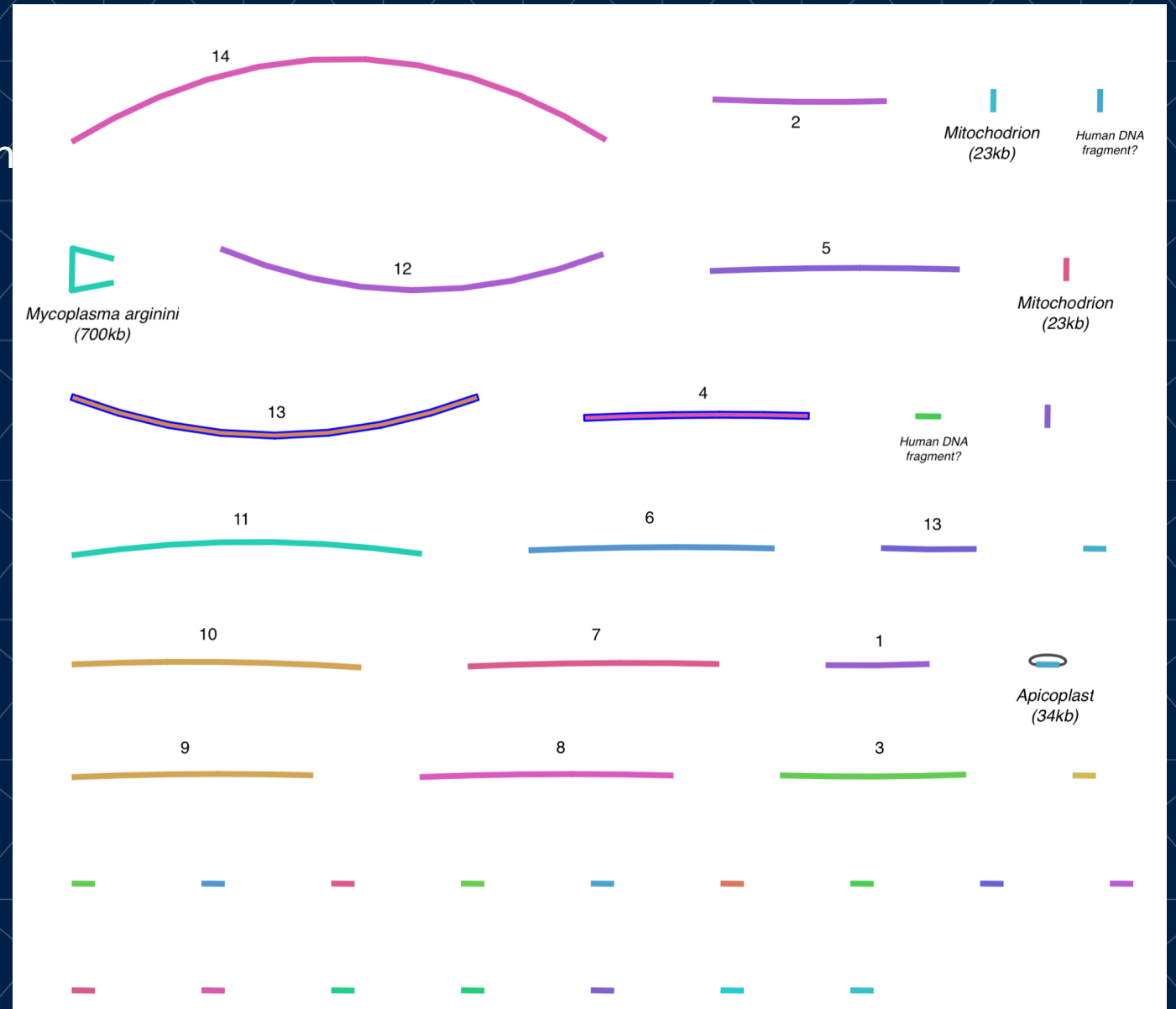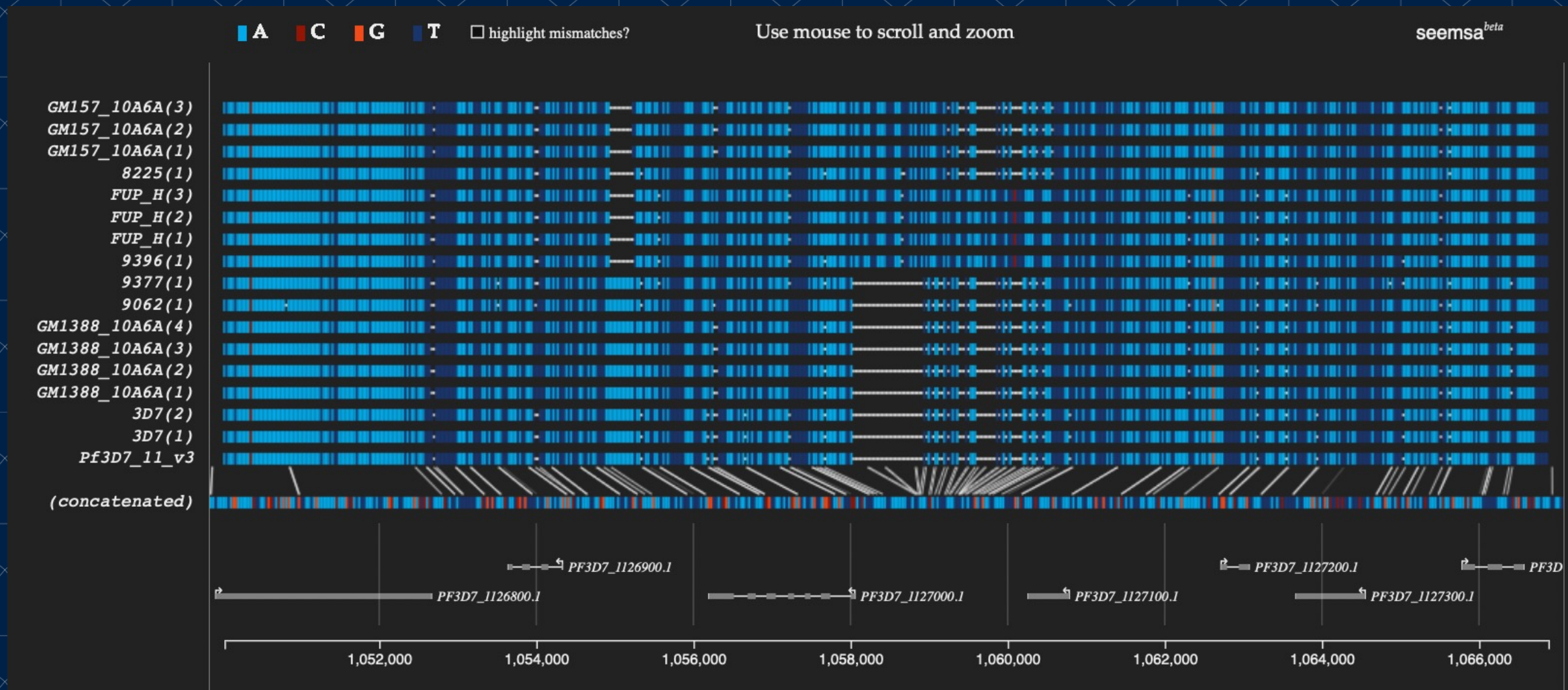2nd attempt: Pacbio whole-genome sequencing

3D7
FUP-H
4 Kenyan parasites
2 Gambian parasites
1 parasite from single-cell sorting

Carried out by James Docker and Amy Trebes, Oxford Genomics Centre for a test of new fragmentation protocol.

Worked amazingly well

Multiple sequence alignment of *P.falciparum* whole genomes

# Acknowledgments

**wellcome centre human genetics**

John Todd
Julian Knight
Andrew Brown
Tony Cutler
Connor Davison
Jia-yuan Zhang
David Smith
Annie Forster
Qijing Shen
Jason Hendry
Hitomi Kuwabara
David Buck
Paolo Piazza
Helen Lockstone

**KEMRI-Wellcome Kilifi, Kenya**
Alexander W. Macharia
Patrick Ombati
Silvia Kariuki

CIMR:
Julian Rayner

OGC:
Amy Trebes
James Docker
David Buck

**PacBio**

Riki Aydeniz
Eirini Maria Lampraki
Mike Eberle
Cillian Nolan

…and HV31.

Dominic Kwiatkowski
1953-2023

**Oxford NANOPORE Technologies**

Simon Mayes
Philipp Reschender
Tonya McSherry

**MalariaGEN**

Dominic Kwiatkowski
Ellen Leffler
Kirk Rockett

*"We thank the patients and staff at the Paediatric Department of the Royal Victoria Hospital in Banjul, Gambia, and at Kilifi County Hospital and the KEMRI-Wellcome Trust Research Programme, Kilifi, Kenya for their help with this study."*

**MalariaGEN**
GENOMIC EPIDEMIOLOGY NETWORK

**Thanks!**