# Characterizing and Correcting Errors in Long Read Sequencing Data

PRESENTER: **Jonn Smith**

## INTRO

The PacBio HiFi protocol can generate consensus reads of quality Q20-Q30 by redundant sequencing of circularized fragments. However, this consensus process is only successful if the HiFi adapter is correctly isolated and removed, and if the circular sequencing completes at least two passes of the original input molecule. Reads failing these criteria are not corrected. Empirically this amounts to, on average, 50% the reads being rejected from the HiFi consensus process and left at their original lower quality.

We use the corrected HiFi reads to construct a graph based on their alignment to the reference. We then align the rejected reads to this graph, favoring edges from the HiFi data and omitting poorly supported edges in the rejected data. This effectively corrects the rejected reads to the same approximate quality as the corrected HiFi reads using the these HiFi reads as a prior.

We believe this method is applicable to other consensus-based sequencing techniques as well (such as R2C2 as applied to Oxford Nanopore data).
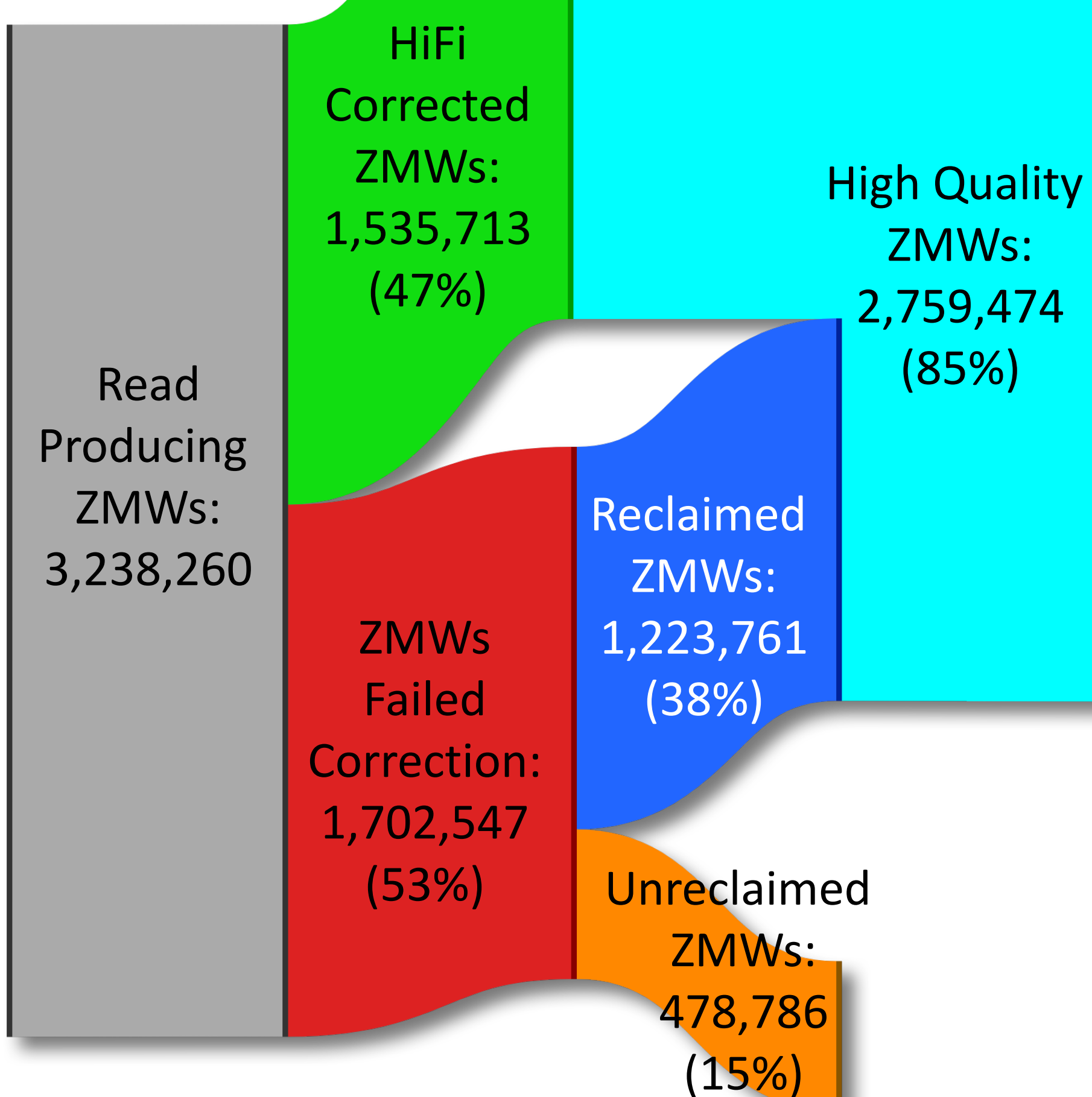
Here we present our prototype of this reclamation process.

## METHODS

1. Create *Partial Order Alignment Graph* from HiFi reads aligned to the reference
2. Align rejected reads to HiFi read graph
3. Correct and extract rejected read sequences based on alignment with graph
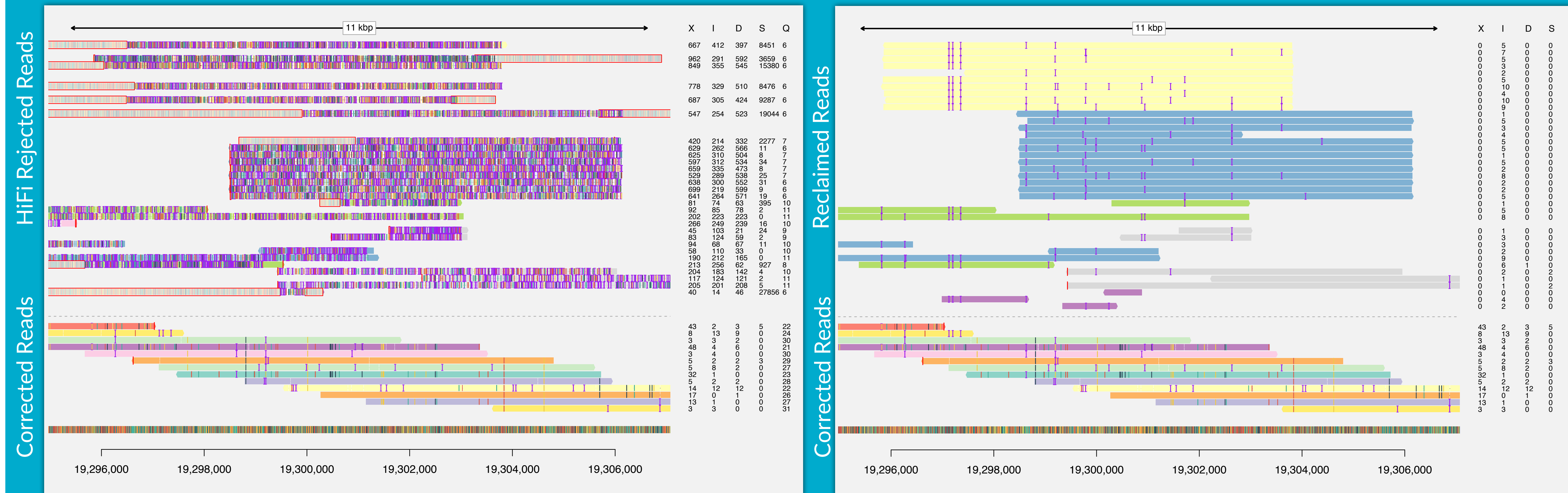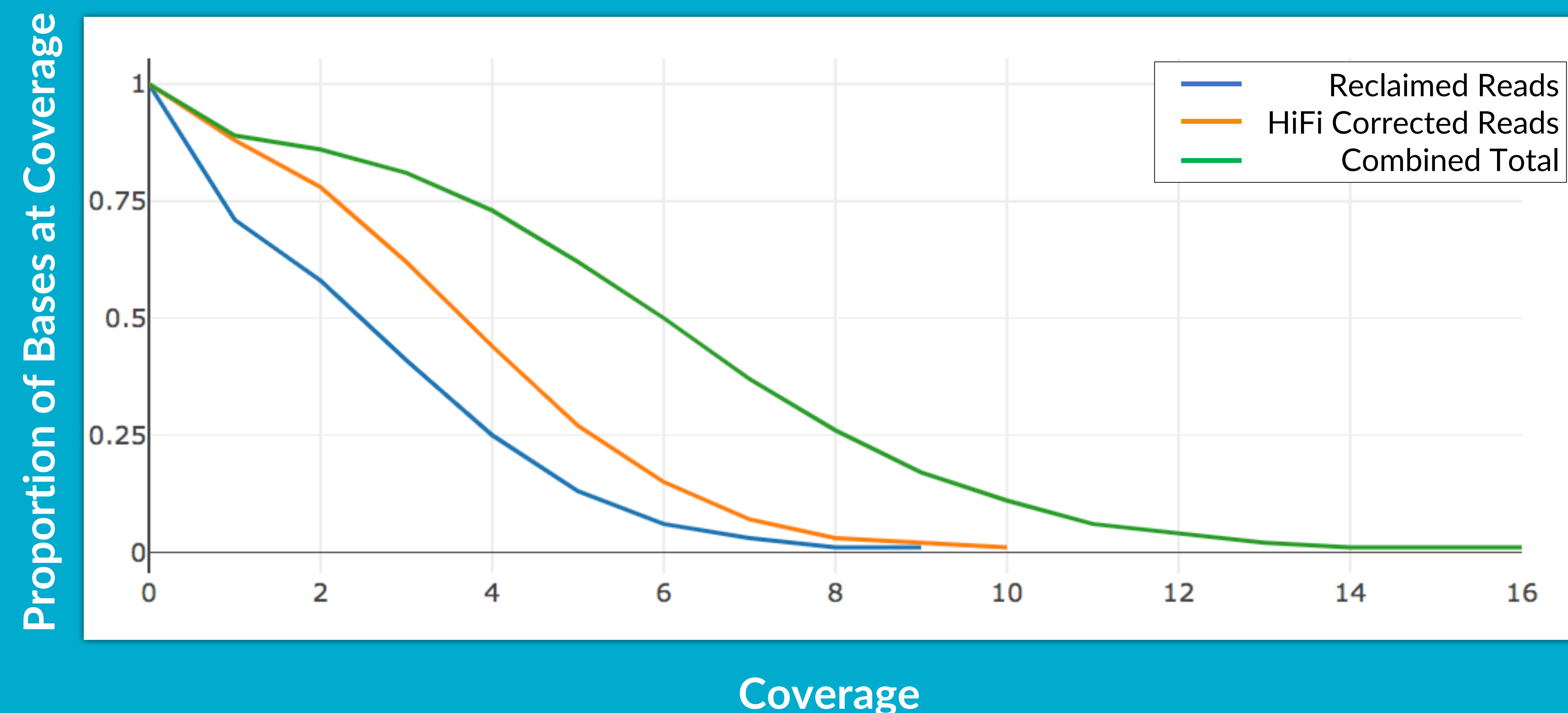4. Select "best" reclaimed read from each ZMW

## RESULTS

### NA19240

Computational Biology and Data Science

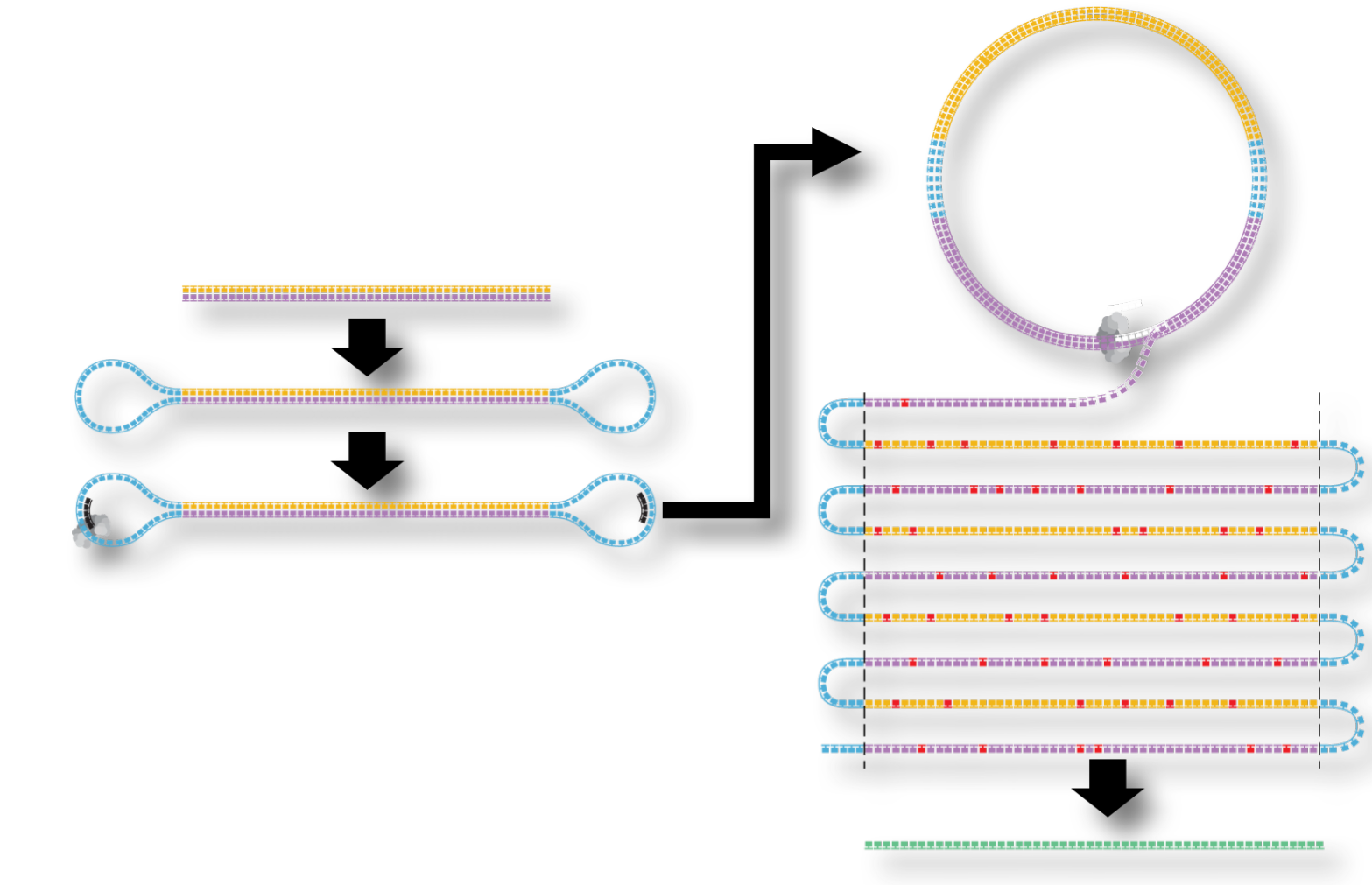Graph-based read correction reclaims 40-90% of reads originally rejected from PacBio HiFi data without resequencing.

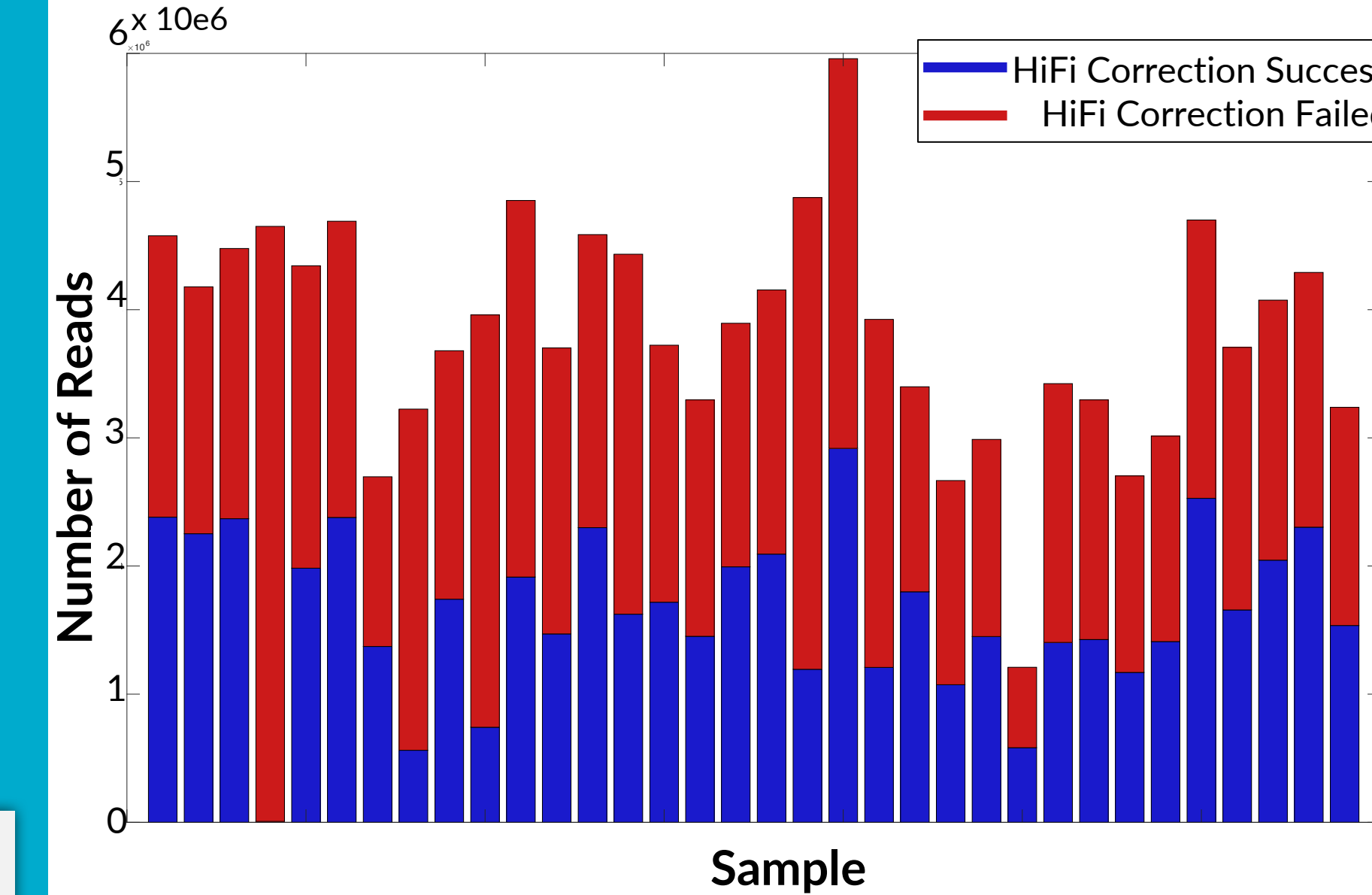X – Mismatches, I – Insertions, D- Deletions S – Soft Clips, Q – Quality
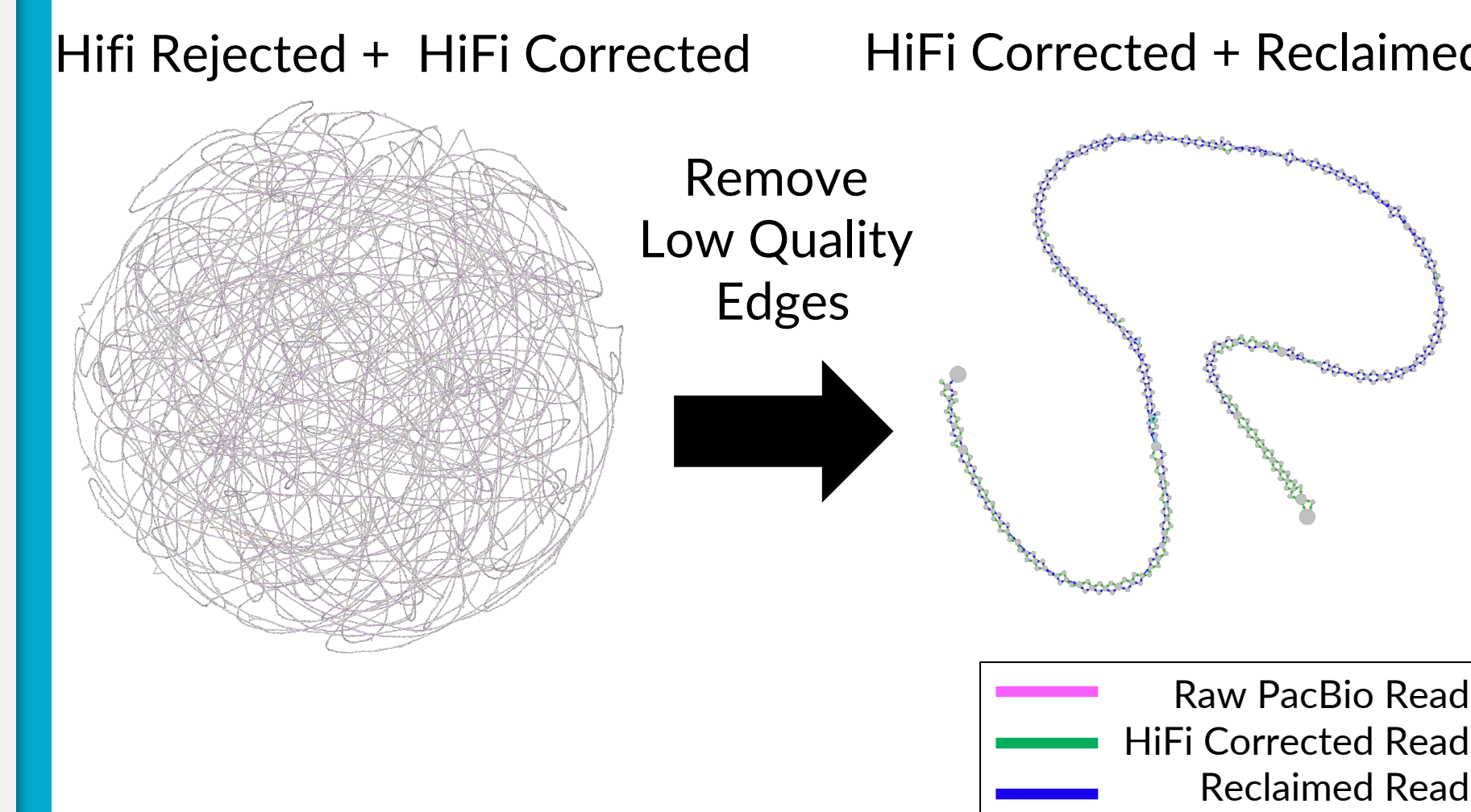

NA19240 Coverage

### PacBio HiFi Sequencing Process



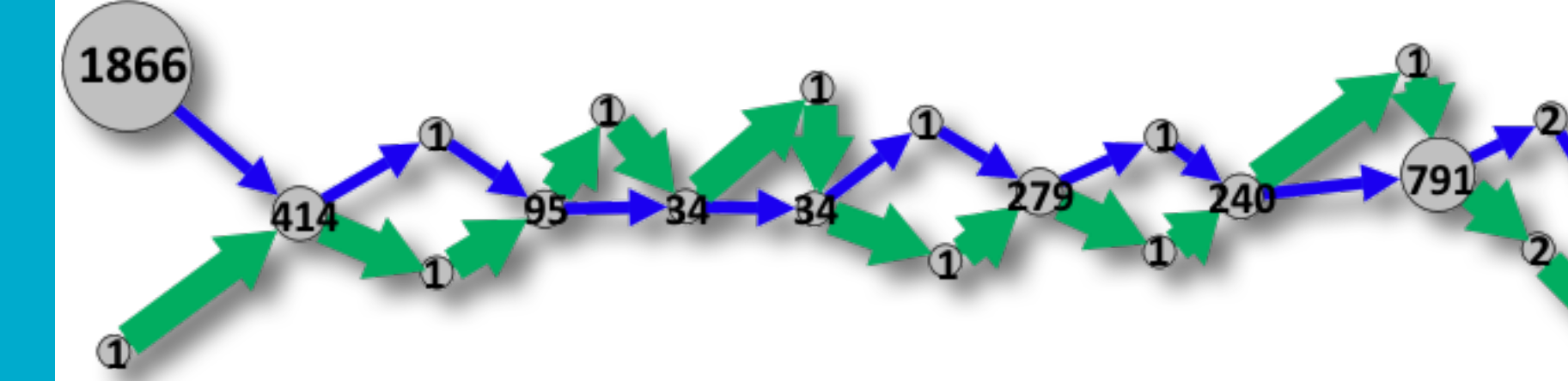### HiFi Read Correction Counts



### Partial Order Alignment Graphs of Reads



Hifi Rejected + HiFi Corrected → HiFi Corrected + Reclaimed

Remove Low Quality Edges

Raw PacBio Reads
HiFi Corrected Reads
Reclaimed Reads

HiFi Corrected + Reclaimed (Graph Start)

**Jonn Smith**, Michael Gatzen, Steve Huang, Maura Costello, Tera Bowers, Kiran Garimella

**BROAD INSTITUTE**

**References:**
- Lee, C. (2003). Generating consensus sequences from partial order multiple sequence alignment graphs. *Bioinformatics*, *19*(8), 999–1008. https://doi.org/10.1093/bioinformatics/btg109
- Rautiainen, M., & Marschall, T. (2019). Graphaligner: Rapid and versatile sequence-to-graph alignment. *BioRxiv*, 810812. https://doi.org/10.1101/810812
- Wenger, A. M., Peluso, P., Rowell, W. J., Chang, P.-C., Hall, R. J., Concepcion, G. T., ... Hunkapiller, M. W. (2019). Highly-accurate long-read sequencing improves variant detection and assembly of a human genome. *BioRxiv*, 519025. https://doi.org/10.1101/519025