# Efficient de novo Assembly of Eleven Human Genomes Using PromethION Sequencing and a Novel Nanopore Toolkit
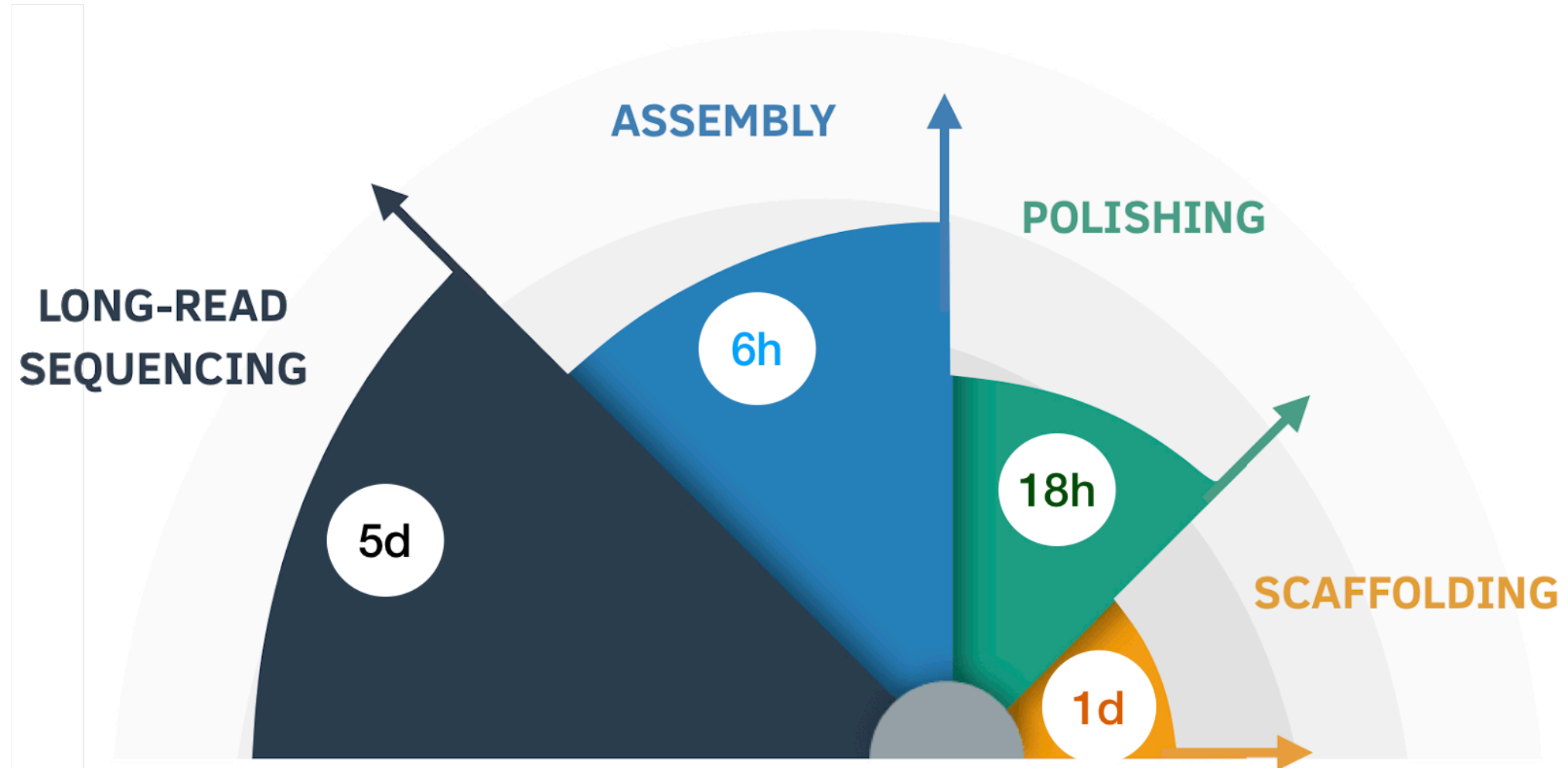
**Miten Jain, Hugh Olsen**

**UC Santa Cruz**
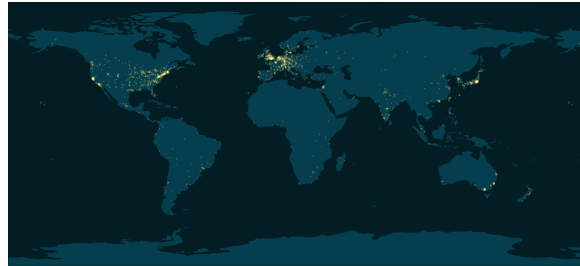
**1 March 2020**

@mitenjain

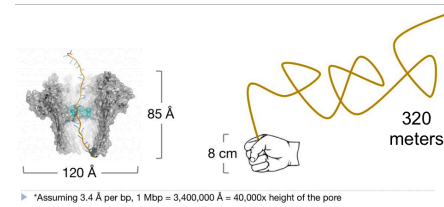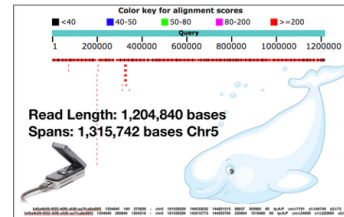# MESSAGE: Reference-quality human genome in ~7 days using nanopore + HiC

# Genomics Using Nanopore Devices

1. Mobile

2. Long read lengths

   Read Length: 1,204,840 bases
   Spans: 1,315,742 bases Chr5

   85 Å
   120 Å
   8 cm
   320 meters

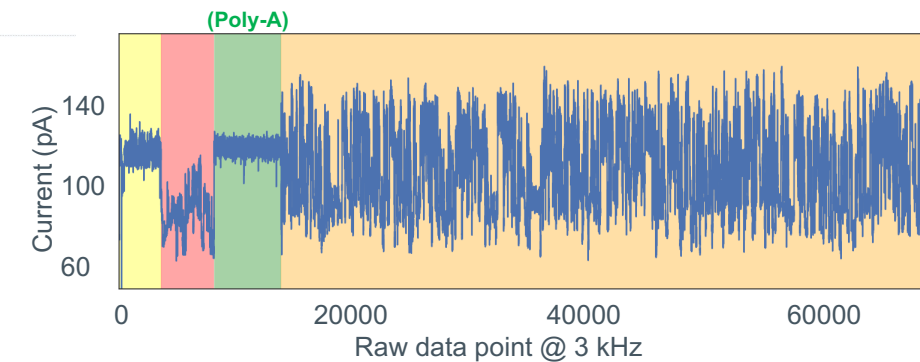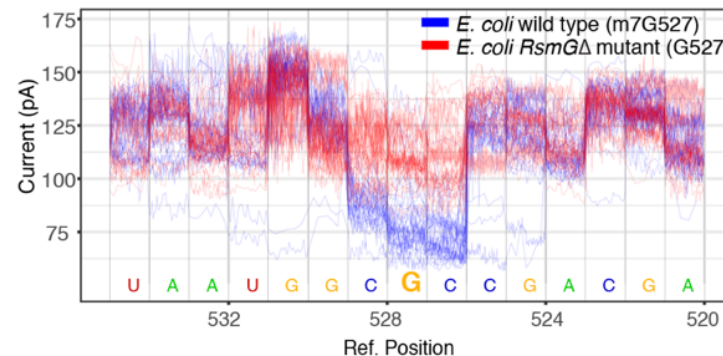   *Assuming 3.4 Å per bp, 1 Mbp = 3,400,000 Å = 40,000x height of the pore

3. Direct readouts of RNA

   (Poly-A)

4. RNA/DNA base modification detection

# What's the bottleneck in capturing genome variation?

- Need for hundreds of high-quality reference genomes

- Sequencing cost

- Sequencing speed

- Scalable and cheaper informatics

# Solution

- Nanopore 100kb+ sequencing

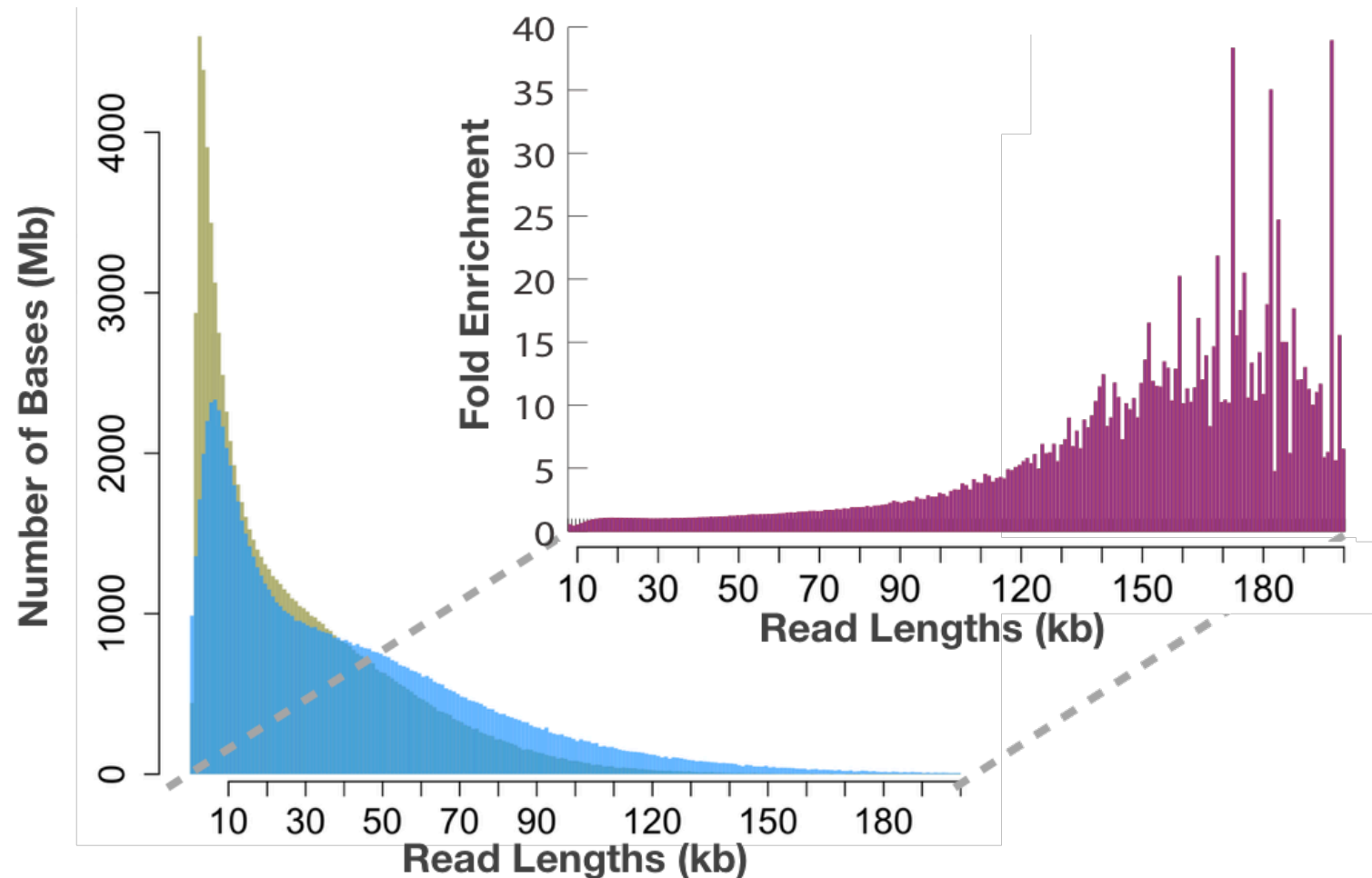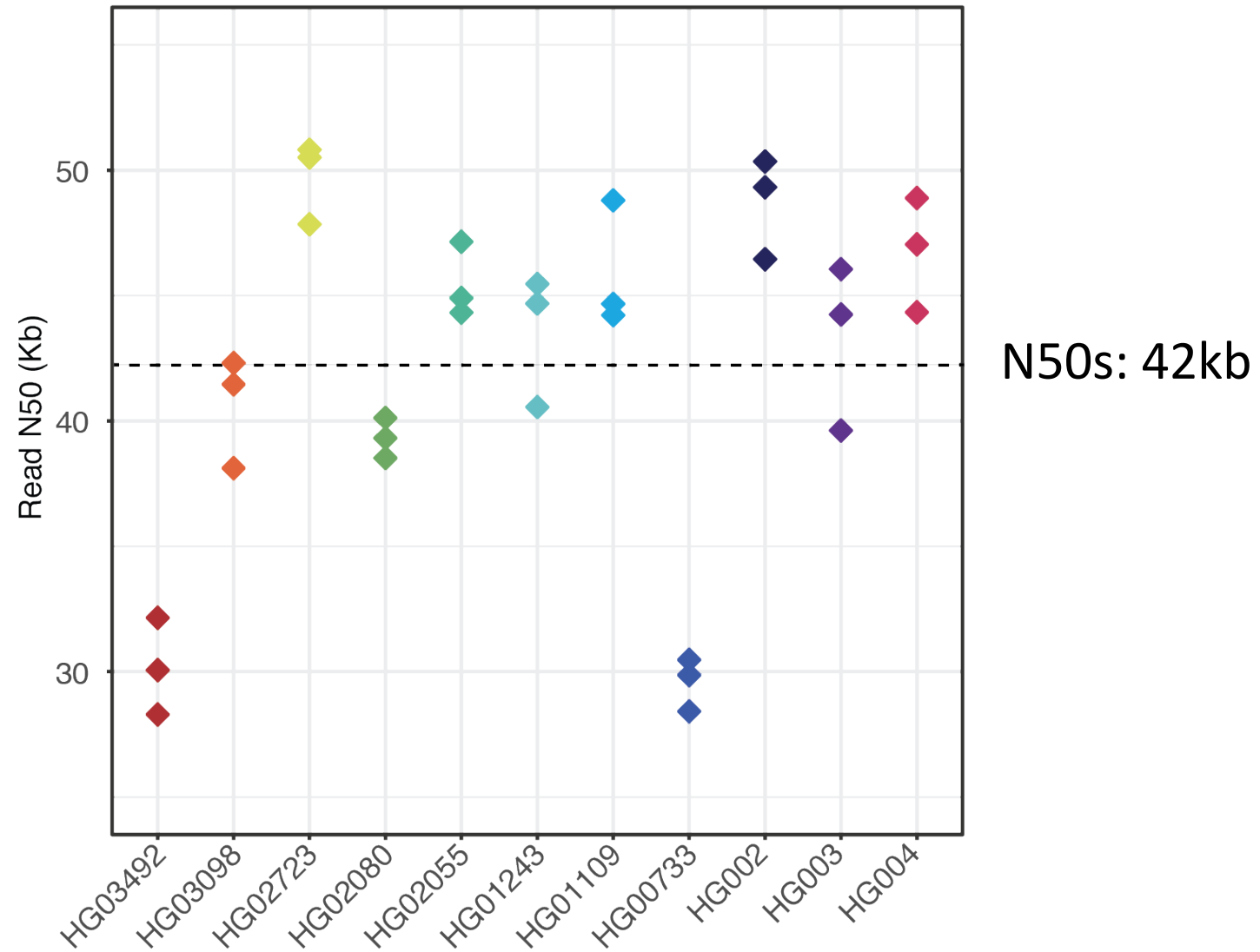- Scalable assembly and polishing

  tools

# Nanopore 100kb+ sequencing

Data acquisition for 11
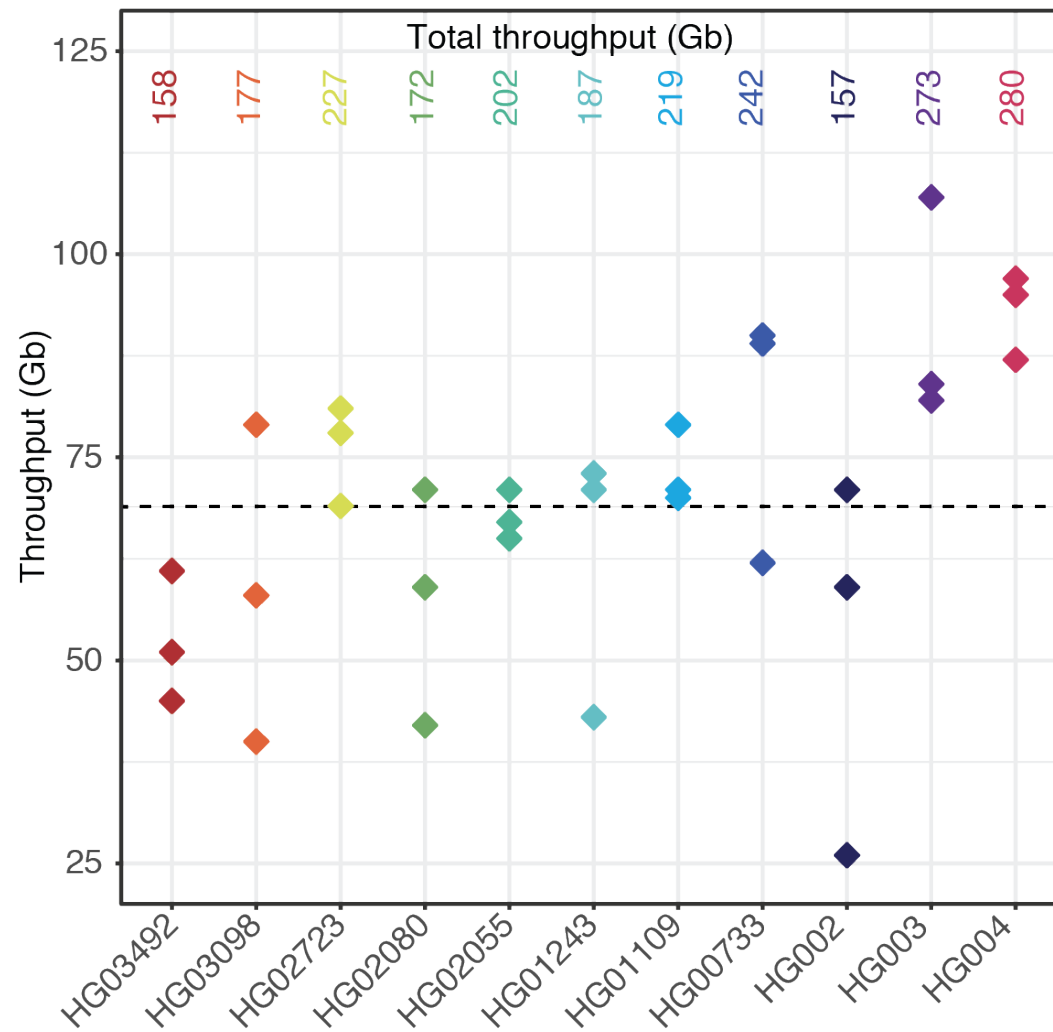genomes in 9 days (>60x
total coverage)

# 7x enrichment of reads >100kb using Circulomics SRE
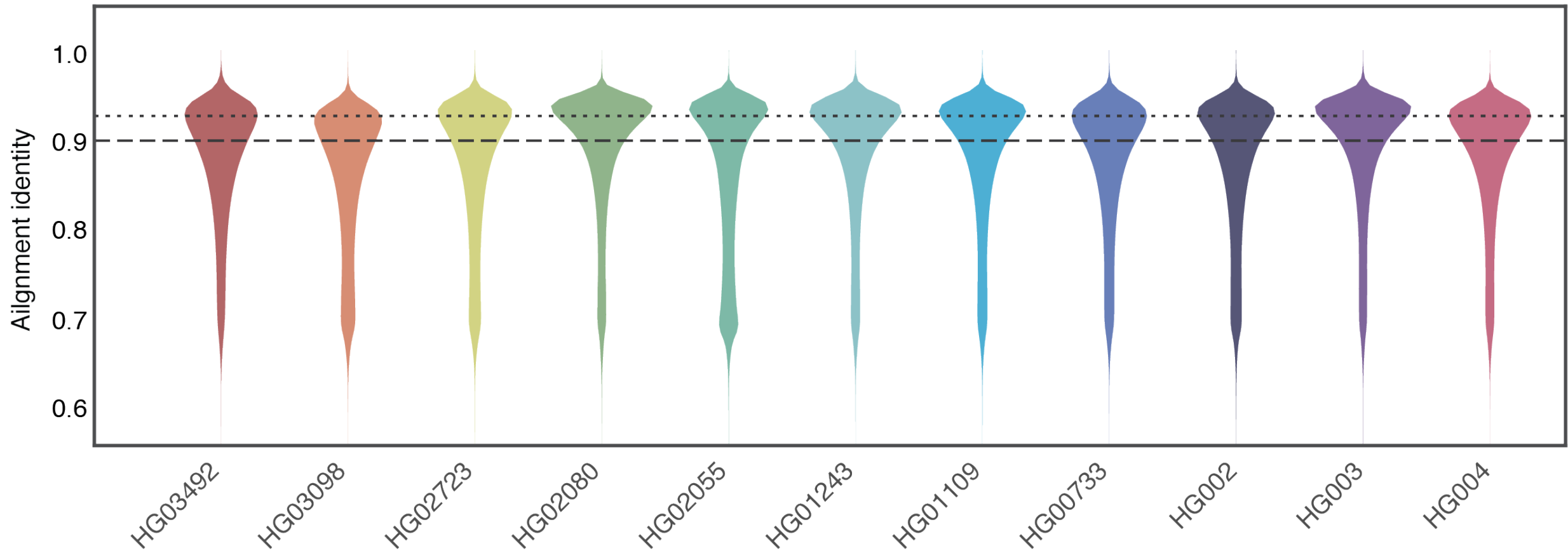


Short Read Eliminator Kit (https://www.circulomics.com)

# Read N50 improvement is reproducible



N50s: 42kb

# PromethION sequencing throughput

# Median alignment identity is 90%

# Scalable assembly and polishing tools

# Pipeline



**Sequencing/Basecalling** → Guppy Flip-Flop

**Assembly** → Shasta

**Polishing** → Margin-Polish → HELEN

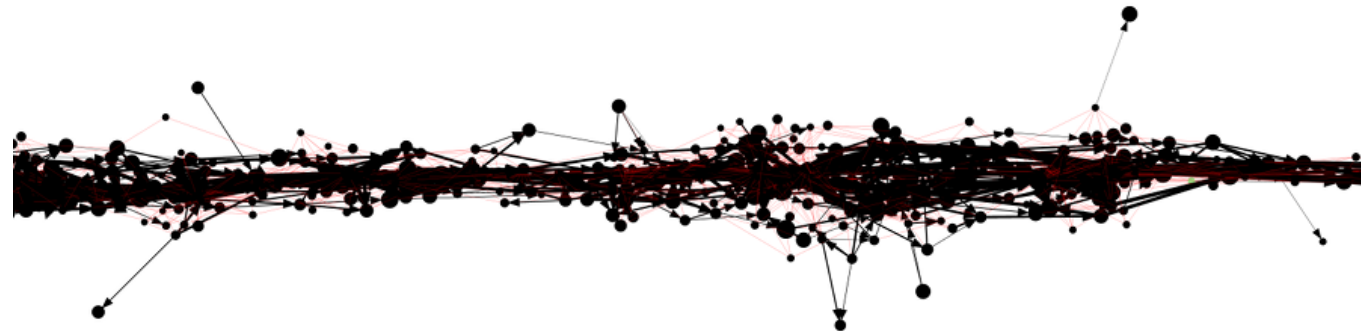**Scaffolding** → HiRise

**Phasing** → FINISHED ASSEMBLY

amazon web services™

# Shasta – a nanopore *de novo* long read assembler
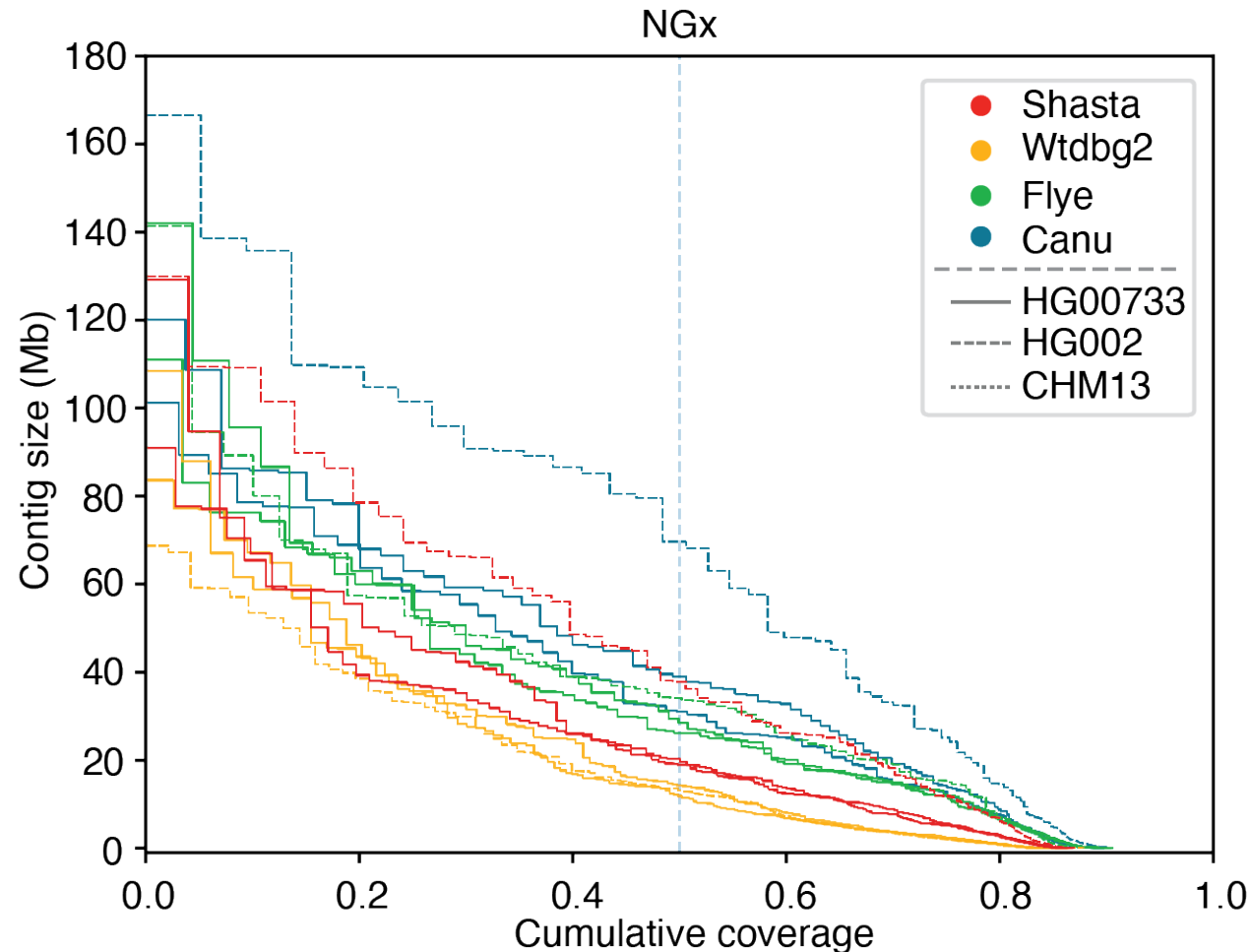
- Extremely fast (can assemble a human genome in < 6 hours on a single node)

- Uses run-length representation of read sequence

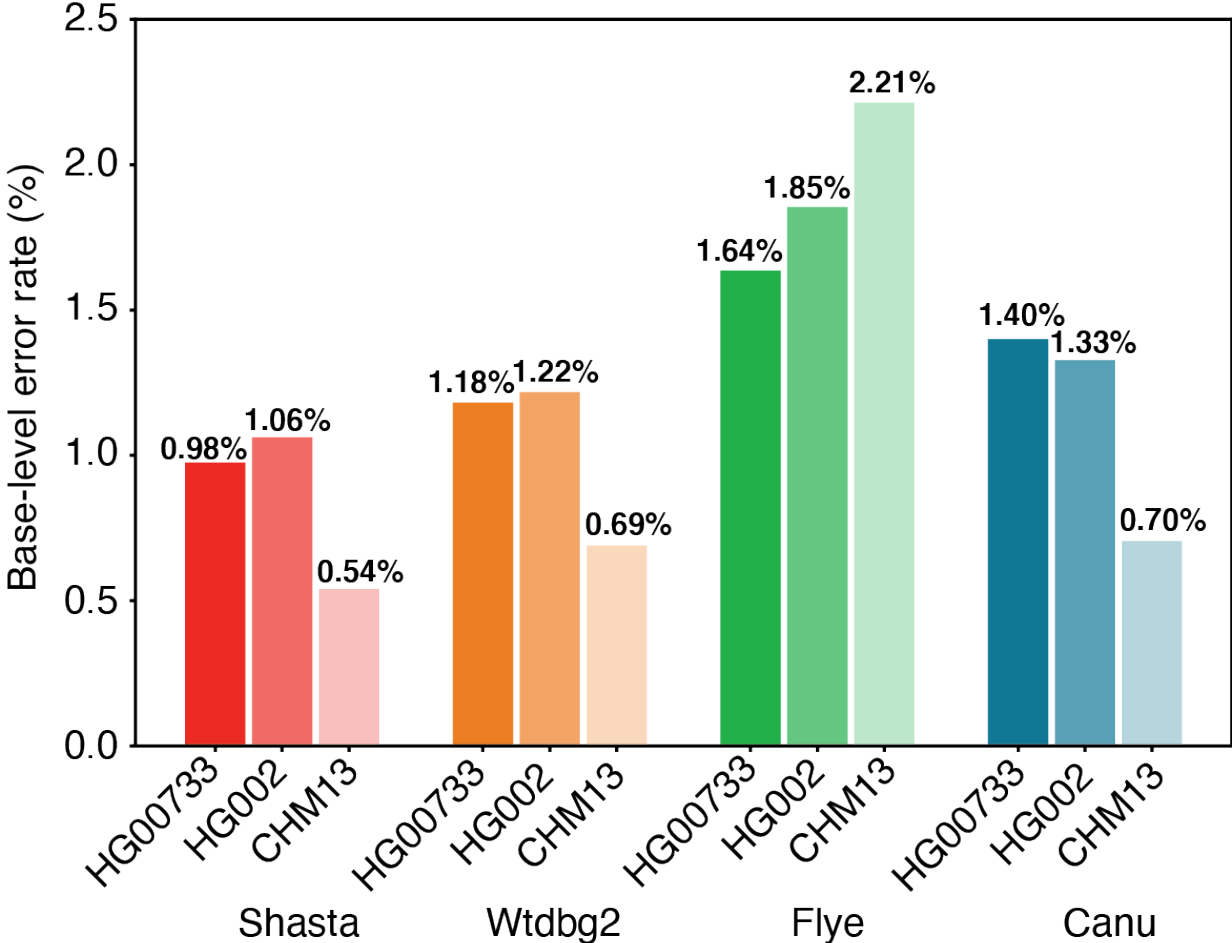- Uses a memory-based marker representation for fast compute

- https://github.com/chanzuckerberg/shasta

# Shasta assemblies are reproducible, with comparable contig NG50



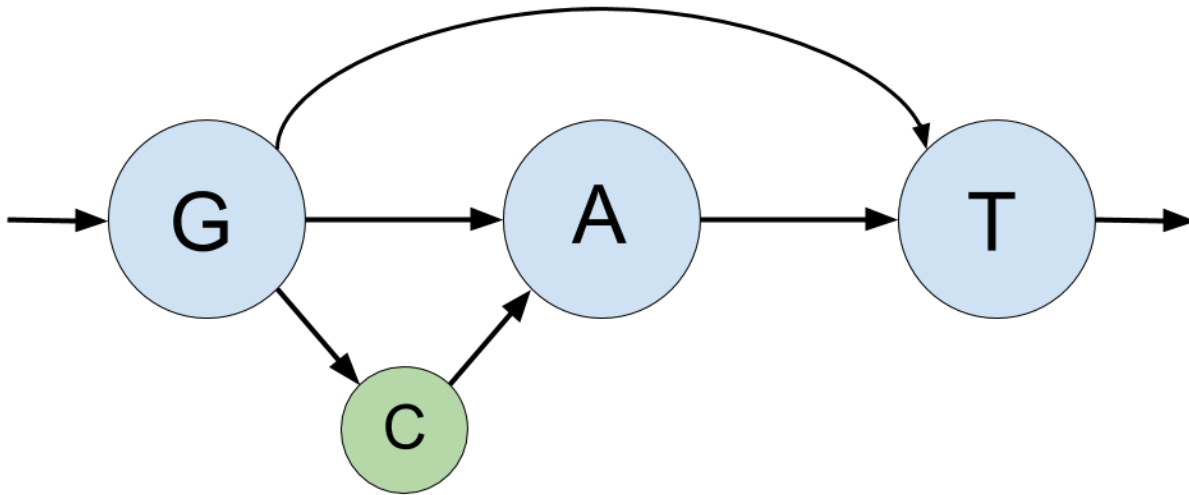Median contig NG50 = 23 Mb

# Shasta assemblies have higher accuracy

# Assembly at a fraction of time and cost



Average cost ($)

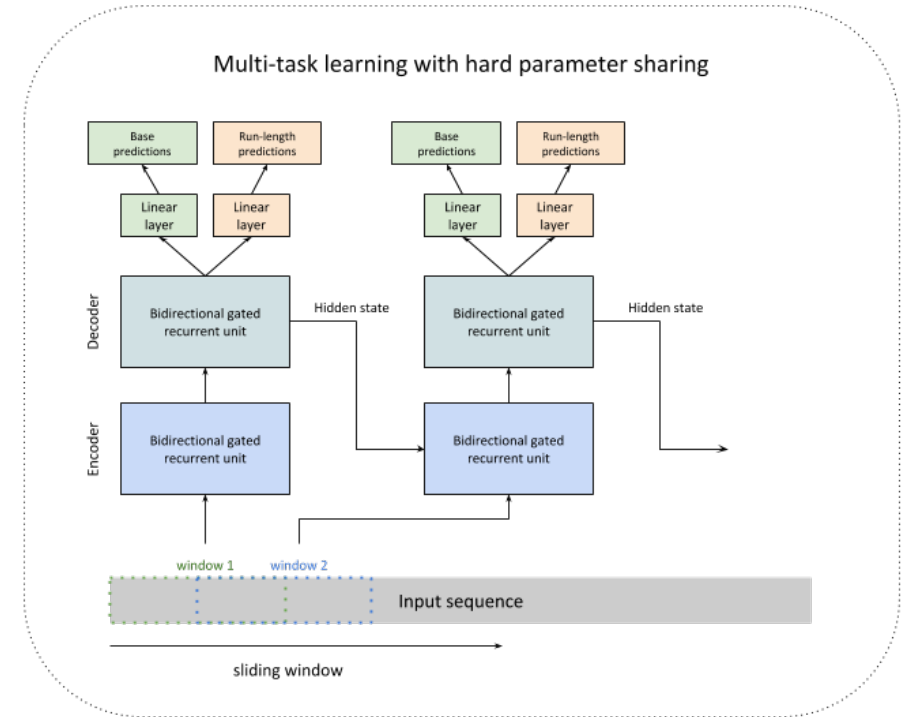| | | | | | |
|---|---|---|---|---|---|
| 0 | 160 | 320 | 480 | 640 | 800 |

Flye — $695.21 — 62.53 hr

Wtdbg2 — $142.03 — 39.28 hr

Shasta — $70 — 5.25 hr

Average elapsed runtime (hour)

16

# Two-step polishing of assemblies

## 1. MarginPolish



A graph-based alignment polisher

https://github.com/UCSC-nanopore-cgl/marginPolish

## 2. HELEN



A DNN-based consensus sequence polisher

https://github.com/kishwarshafin/helen

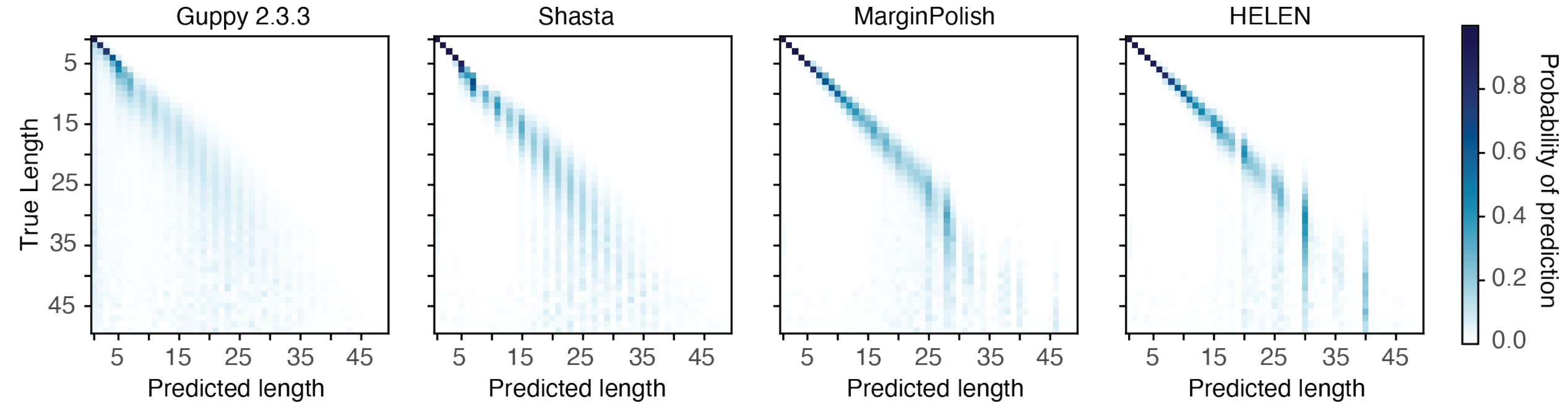UC SANTA CRUZ

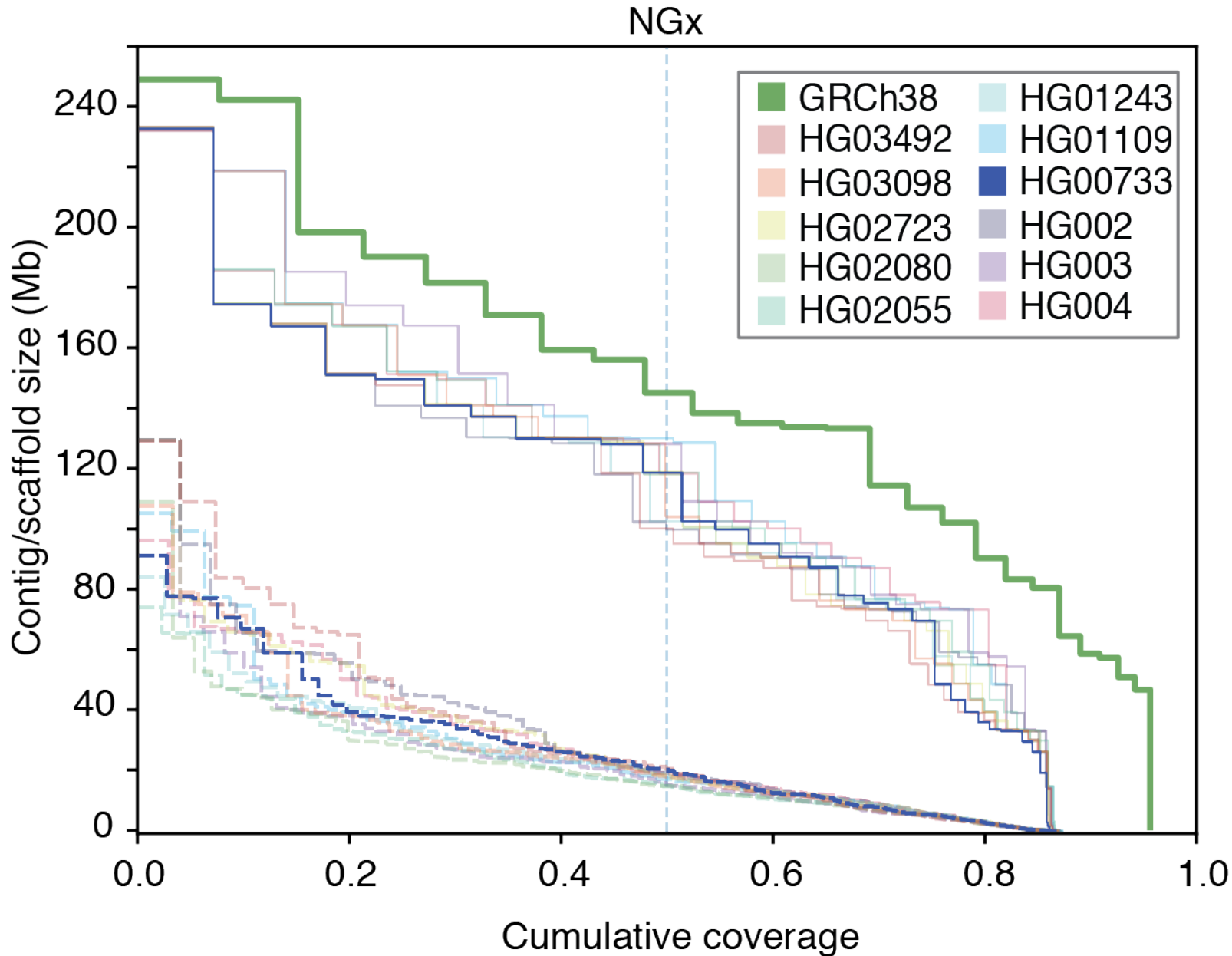# MarginPolish and HELEN outperform other polishers

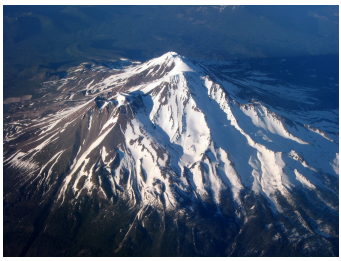# Polishing at a fraction of time and cost

# Improvements in homopolymer length predictions

# Chromosome-level scaffolding using HiC data
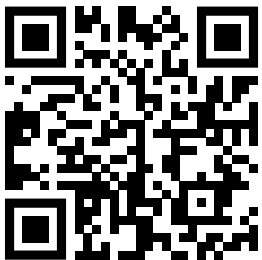
# Key next steps

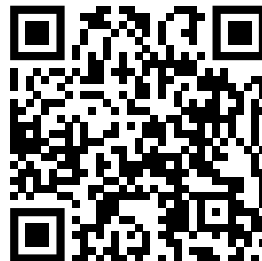- Faster basecalling (ONT)

- Haplotype phasing (UCSC, CZI)

Shasta

MarginPolish

HELEN

# Data resources

- Human Pangenome Reference Consortium
    - https://tinyurl.com/hpp-hg002
    - https://tinyurl.com/hpp-data

- Human Pangenomics
    - https://tinyurl.com/hpgp-data

- Nanopore DNA consortium
    - https://tinyurl.com/na12878-dna

- Nanopore RNA consortium
    - https://tinyurl.com/na12878-rna

# Acknowledgements and Collaborators

U British Columbia
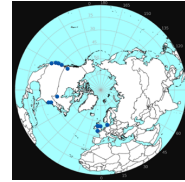
U Oregon

UCSF

UC Santa Cruz

UCLA

Penn

Johns Hopkins

Daniel Garalde
Rosemary Dokos
Simon Mayes
Vania Costa
Jon Pugh
Chris Seymour
Chris Wright
David Stoddart
Dan Turner

National Human Genome Research Institute

amazon web services™

OICR

U Birmingham

U Nottingham

Newcastle U

U Copenhagen

NASA

CHAN ZUCKERBERG INITIATIVE

## UC Santa Cruz

| | | |
|---|---|---|
| Hugh Olsen ✓ | Mark Akeson ✓ | Kishwar Shafin ✓ |
| Robin Abu-Shumays | Karen Miga ✓ | Sofie Salama |
| Vinay Poodari | Benedict Paten ✓ | Marina Haukness |
| Niki Thomas | David Haussler | Trevor Pesout |
| Jenny Vo | Angela Brooks | Ryan Lorig-Roach |
| Logan Mulroney | Manny Ares | Chris Vollmers |
| Ed Green | Jeremy Sanford | Daniel Kim |
| Holger Schmidt | Holger Schmidt | Susan Carpenter |
| Ali Yanik | Ali Yanik | Ed Green |
| | Kristof Tigyi | Ro Kamakaka |

vg

Adam Novak
Glenn Hickey
Jordan Eizenga
Erik Garrison
Jean Monlong
Xian Chang

Adam Phillippy (NHGRI)
Fritz Sedlazeck (Baylor)
Sergey Koren (NHGRI)

circulomics

Kelvin Liu
Duncan Kilburn
Jeffrey Burke

Chan-Zuckerberg

Paolo Carnevali ✓

@mitenjain