Telomere-to-telomere assembly of a complete human X chromosome

Sergey Koren

Staff Scientist, Genome Informatics Section, NHGRI









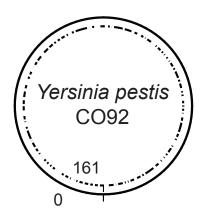
40 years of genome assembly

Rodger Staden (1979)

• "If the 5' end of the sequence from one gel reading is the same as the 3' end of the sequence from another the data is said to overlap. If the overlap is of sufficient length to distinguish it from being a repeat in the sequence the two sequences must be contiguous. The data from the two gel readings can then be joined to form one longer continuous sequence."



- First complete de novo assemblies
 - 2012: Bacteria (10⁶ bp)





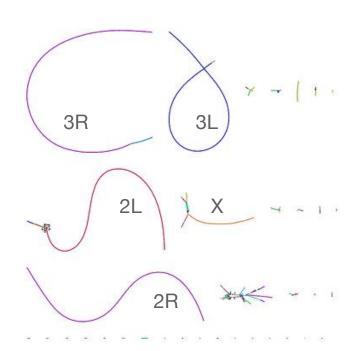
- First complete de novo assemblies
 - 2012: Bacteria (10⁶ bp)
 - 2014: Yeast (10⁷ bp)





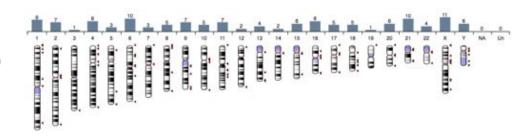
• First complete de novo assemblies

- 2012: Bacteria (10⁶ bp)
- 2014: Yeast (10⁷ bp)
- 2014: Drosophila (10⁸ bp)



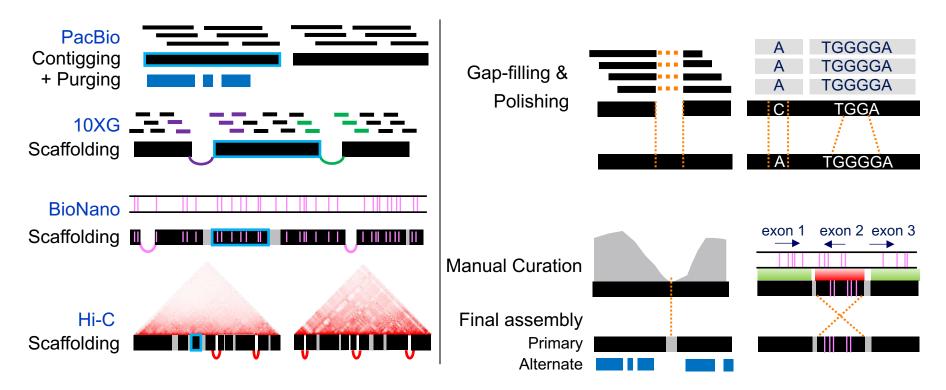


- First complete de novo assemblies
 - 2012: Bacteria (10⁶ bp)
 - 2014: Yeast (10⁷ bp)
 - 2014: Drosophila (10⁸ bp)
 - ????: Human (10⁹ bp)





The Vertebrate Genomes Project Pipeline





VGP GenomeArk: 1st data release





Jennifer Vashon of Maine Department of Inland Fisheries and Wildlife, left, and UMass lynx team coordinator, Tanya Lama, with an adult male lynx from northern Maine whose DNA was used to create first-ever whole genome for the species. The lynx has since been released to the wild. (MassWildlife photo / Bill Byrne)



Amblyraja radiata thorny skate data i taxid



Lynx canadensis Canada lynx data | taxid



Anabas testudineus climbing perch data | taxid



Mastacembelus armatus tire-track eel data I taxid



Vichocentrus centrarchus flier cichlid data I taxid



Ornithorhynchus anatinus platypus data | taxid



Astatotilapia calliptera eastern happy data I taxid



Phyllostomus discolor pale spear-nosed bat data I taxid



Calypte anna
Anna's hummingbird
data I taxid



Rhinatrema bivittatum two-lined caecilian data I taxid



Cottoperca gobio
Channel bull blenny
data I taxid



Rhinolophus ferrumequinum greater horseshoe bat data | taxid



Goode's thornscrub tortoise data I taxid



Strigops habroptilus kakapo data I taxid



Gouania willdenowi blunt-snouted clingfish data | taxid



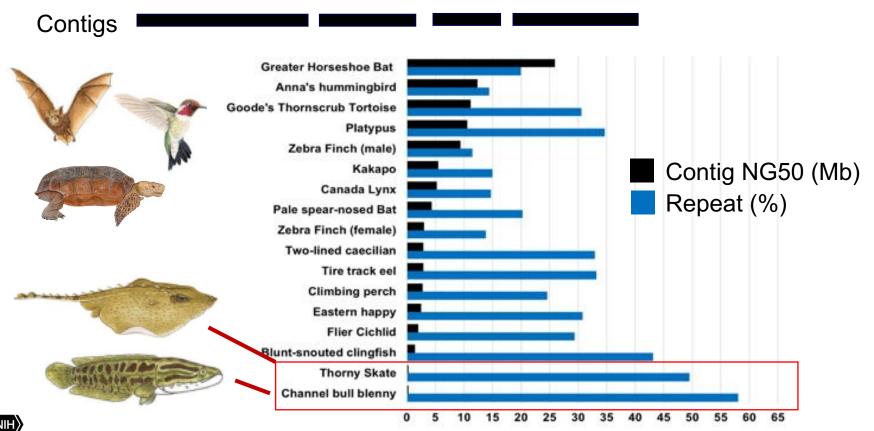
Taeniopygia guttata zebra finch data I taxid



Not all genomes are equal

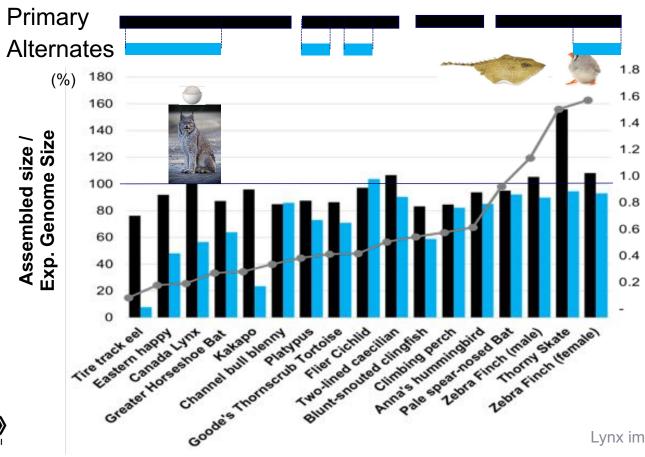
NHGRI





Heterozygosity causes allelic duplication







1.6% Image credit: wikimedia

Falcon- unzip	Primary	Alt
Size (Gbp)	1.95	0.73
NG50 (Mbp)	2.6	0.02
BUSCO	94.2%	40.6%
Duplication	20.8%	3.4%



Resolving haplotypes

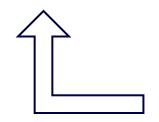


The genomes assembly problem





Duke, highland sire







~1% heterozygosity



Molly, yak dam





State of the art: pseudo-haplotype

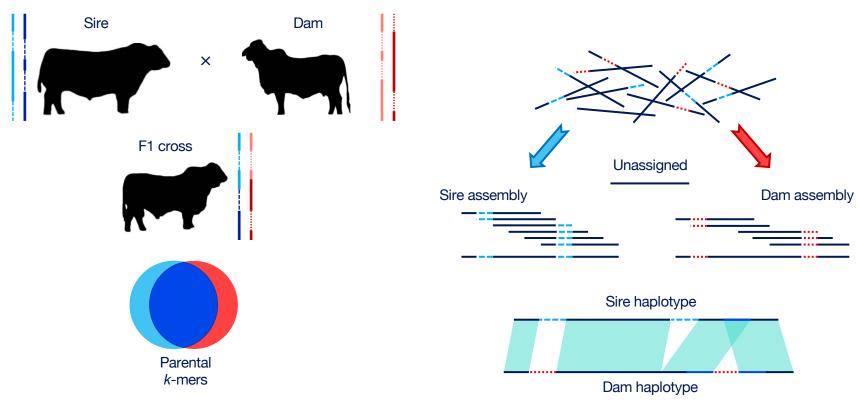






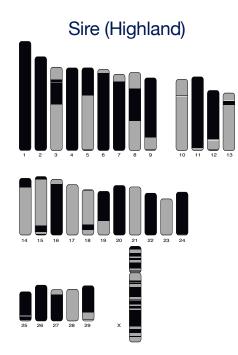
Trio binning with TrioCanu



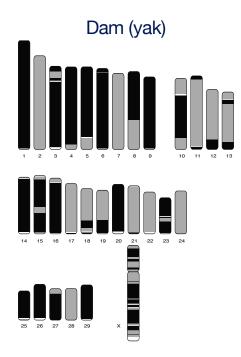




Esperanza: The nearly perfect diploid









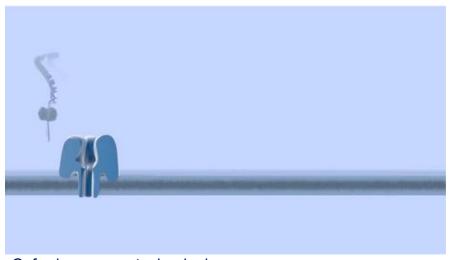
Resolving repeats



Real-time data from a nanopore

\$1000 (free) instrument \$100 / bacterial genome 85–90% read accuracy





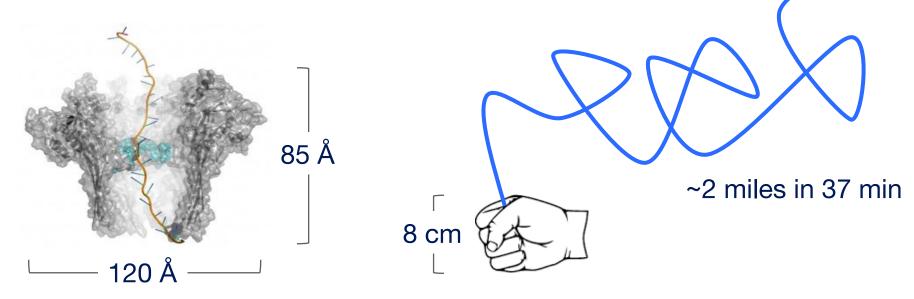
Oxford nanopore technologies



Ultra-long read nanopore sequencing

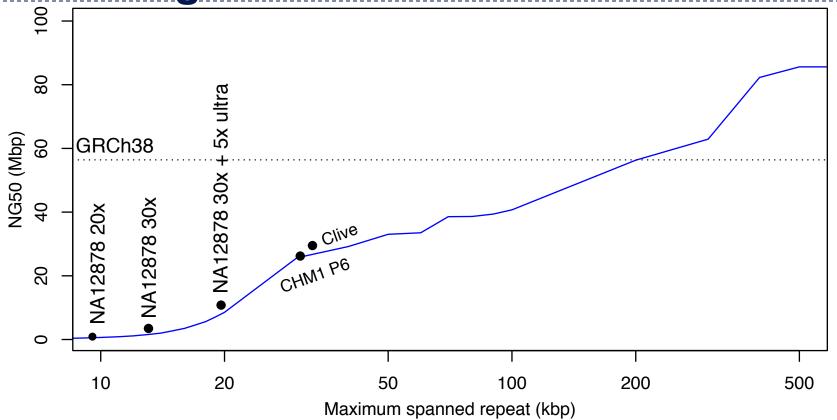
• ONT R9 pore: E. coli CsgG membrane protein

Read lengths >1 Mbp possible





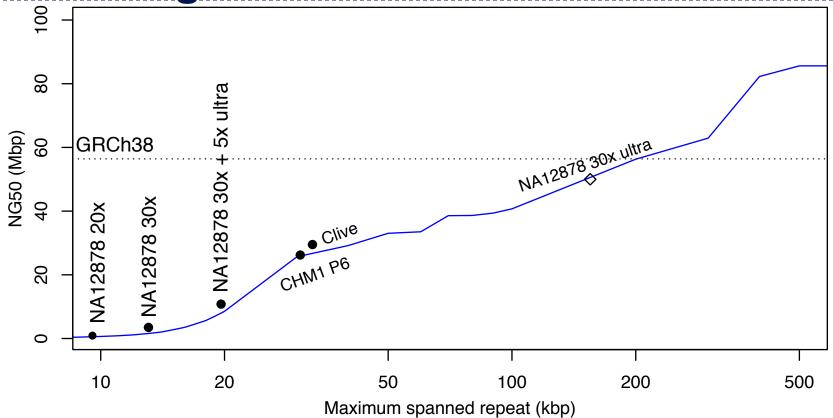
Ultra-long read benefits





Nanopore sequencing and assembly of a human genome with ultra-long reads. Jain, Koren, Miga, Quick, Rand, Sasani, Tyson, et al. *Nature Biotech* (2018)

Ultra-long read benefits





Nanopore sequencing and assembly of a human genome with ultra-long reads. Jain, Koren, Miga, Quick, Rand, Sasani, Tyson, et al. *Nature Biotech* (2018)

The One Genome Project



The human reference is incomplete

- 368 unresolved issues, 102 gaps
- Segmental duplications, rDNAs
- · Centromeres, telomeres, heterochromatin

These gaps contain important information

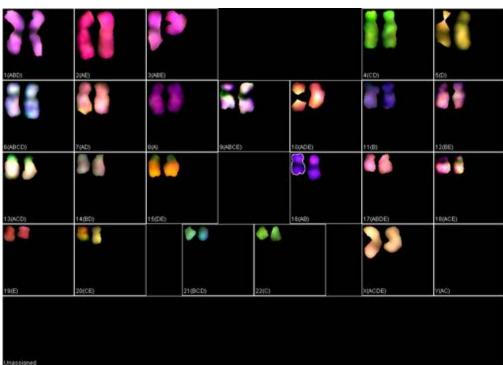
- Missing reference sequence leads to mapping artifacts
- Variation in these gaps is unexplored (e.g. rDNAs)
- We don't know what we don't know...

Long-read sequencing can close these gaps



The One Genome: CHM13hTERT







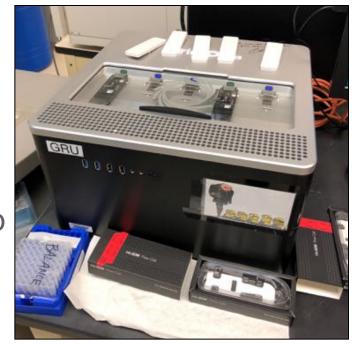
CHM13 sequencing at NISC





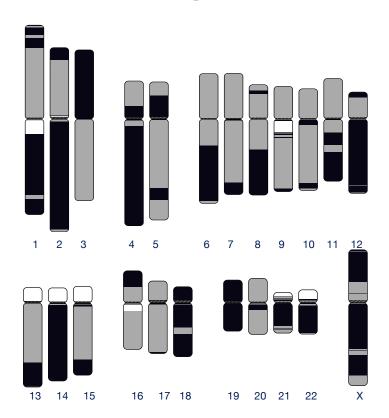


- From May 1/18 Jan 8/19
 - 94 MinION/GridION flow cells
 - 11.1M reads
 - 155 Gb (1.6 Gb / flow cell) (50x)
 - 99 Gb in reads >50 kb (32x)
 - 78 Gb in reads >70 kb (25x)
 - Max mapped read length 1.04 Mb
- Assembled with Canu





The human genome, 2017



GRCh38

The Genome Reference Consortium consists of:





The McDonnell Genome Institute at Washington University



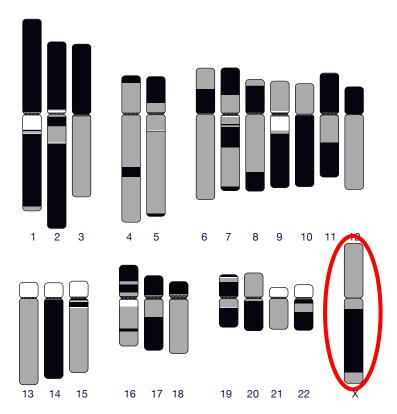
The European Bioinformatics Institute

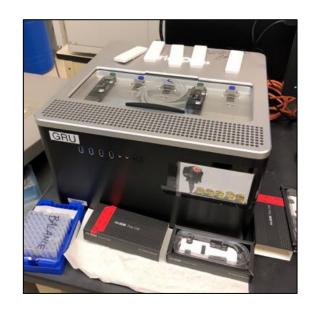


The National Center for Biotechnology Information



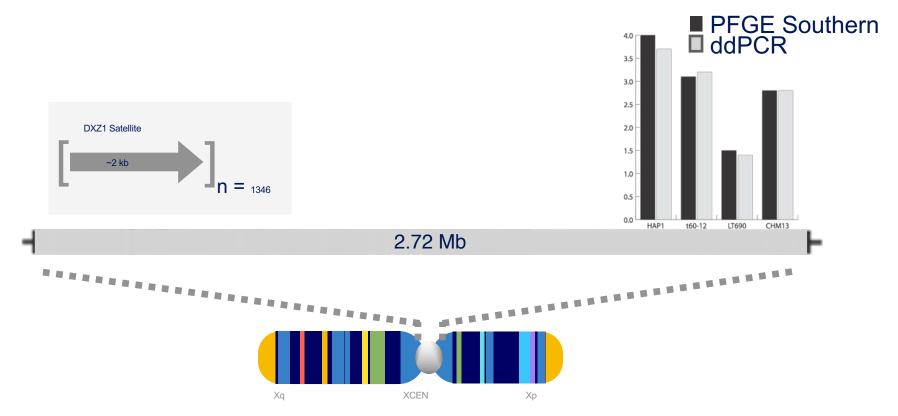
The human genome, 2018





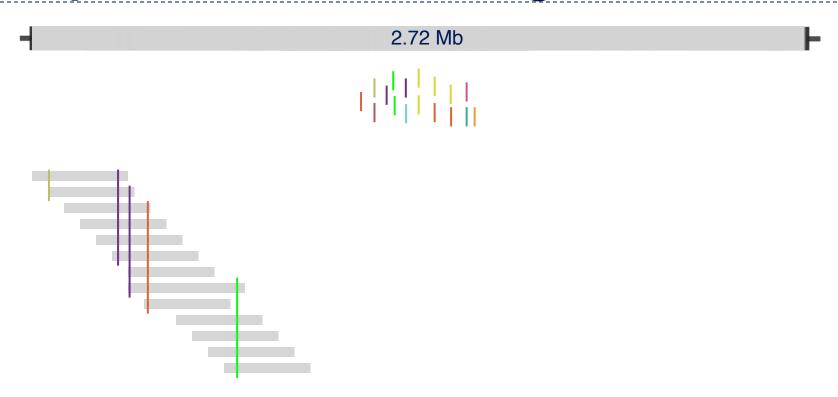


ChrX centromere detail





Repeat resolution using variants





First complete X chromosome

Xq27.1...

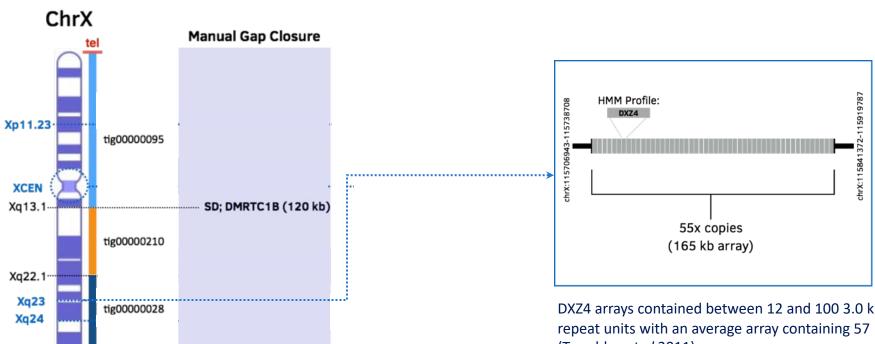
Xq27.3"

tig00002559

tig00344326

SD (134 kb)

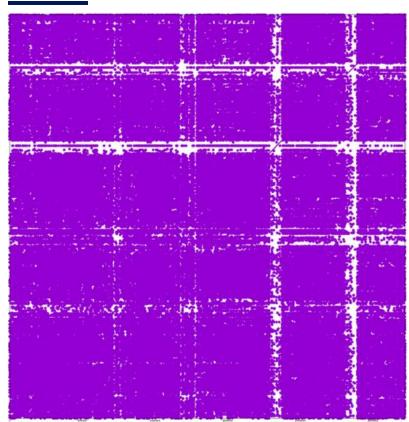




DXZ4 arrays contained between 12 and 100 3.0 kb (Tremblay et al 2011)

How to validate 1300 copy 2kb repeat?

500kb



>=500bp, >=97% idy



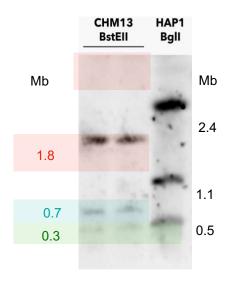
Independent markers confirm structure





Lab validation supports structure







Almost there!

- Large repeats remain challenging
 - Satellite arrays, SegDups, rDNAs, ...
- Haploid assembly is solved by long reads
 - Complete human reference in 1–2 years
- Diploid assembly is solved by trios/Hi-C/ss
 - Complete haplotypes will become the new norm
- Enabling reference genomes for all species
 - Genome 10k and the Vertebrate Genomes Project

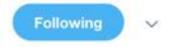


The next 20 years



Bioinformatics is currently artesinal





A useful analogy I used this week:
Sequencing, analysing and interpreting
genomes is 'routine' in the same way the US
Navy 'routinely' lands planes on aircraft
carriers. It might happen regularly by well
trained crew with the right equipment but it is
not an easy thing to do.

3:45 AM - 13 Sep 2018

392 Retweets 974 Likes





Scaling up, training volunteers



Marcela Uliano





Giulio Formenti

Simona Secomandi





Univ. of Milan



Chul Lee



Seoul Nat. Univ.









Maxmilian Driller Calvinna Caswara Majid Vafadar





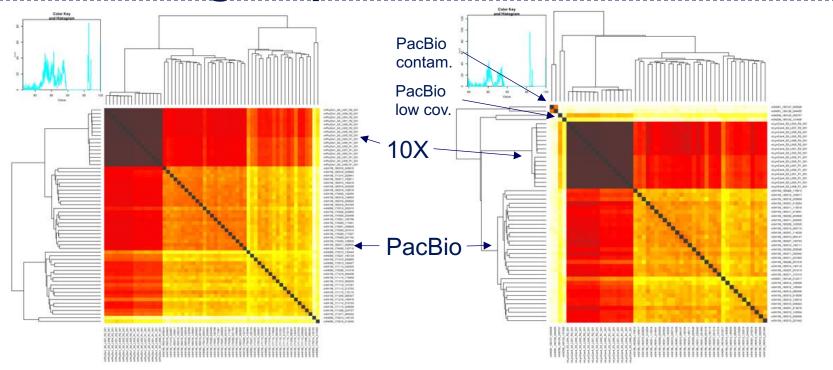




Chai Fungtammasan Nicholas Hill



Automating sequence QC

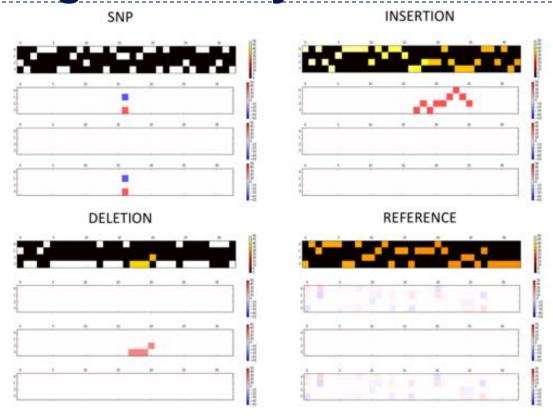


Example of no contamination (mPhyDis1)

Rice SMRT cells in Canadian Lynx



Automating assembly QC?





First trainee results



Carmine Bee-eater

Total bp	1,182.18 Mbp
Number of contigs	723
NG50	10.27 Mbp
Max contig size	41.28 Mbp



Common brushtail possum

Total bp	3,450.80 Mbp
Number of contigs	3,124
NG50	4.60 Mbp
Max contig size	35.90 Mbp



Evolving compute requirements

Whole Genome Applications HPC³

192 physical or 384 hyper-threaded

8 GB

1 TB

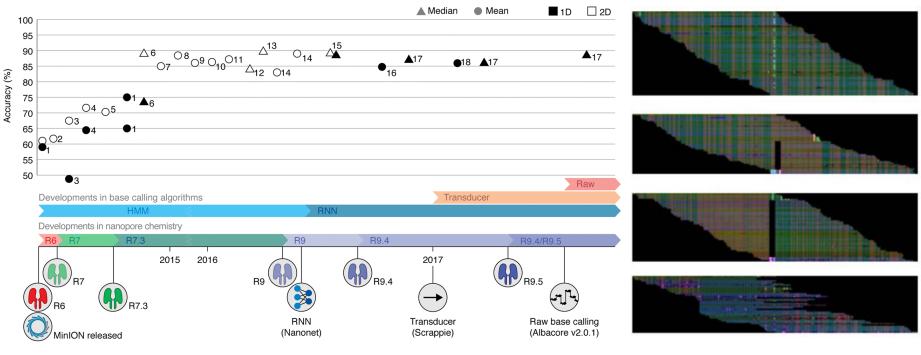
30 TB4 serving 1800 IOPS

70 - 100 TB



These are biowulf norm nodes!

ML gaining popularity

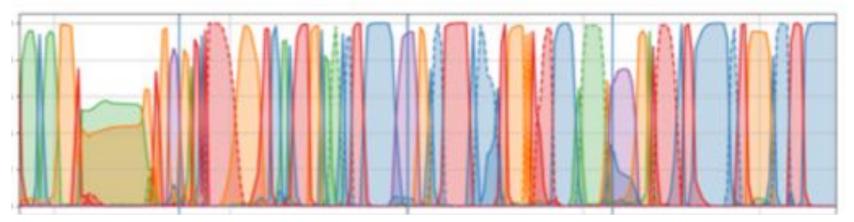


From squiggle to basepair: computational approaches for improving nanopore sequencing read accuracy. Rang et al. 2018





Yet another nanopore basecaller: flappie



- Runtime for 2.5k reads = 15.14h (8 cores)
 - 11m reads / 2.5k = 1/4000th of total
 - 15.14 * 8 * 4000 = 515,282 CPU h!
 - Reservation!

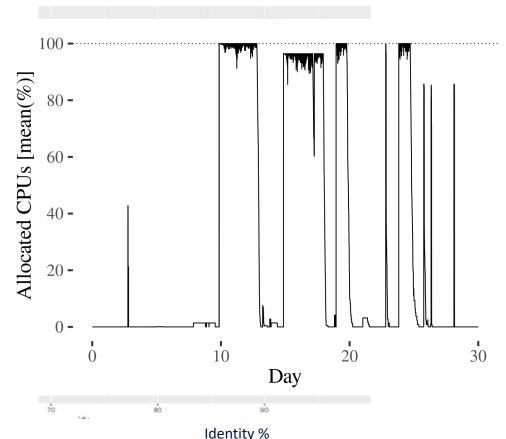


Tools changed too fast

- New GPU basecaller released before reservation
 - Same job took 0.03 h
 - 0.03 * 1 * 4000 = 130 CPU h
 - Finished in < 3 hrs

 Used reservation for other analysis but under-utilized

But at least it is better





Democratized sequencing is here

- \$750 40x human genome*
 - ▶ 15–90m sample prep
 - ▶ 1m–48h runtime
 - 0.1Tb per hour (6Tb in 48h)
 - ▶ 1.44GB/s raw data stream
- >100kb reads possible
- Pay for runs, not instrument
 - "Sequence until"
 - "Read until"
- Primary data to retain?
 - ▶ 10tb/genome @ 10k genomes
 - ▶ 100 petabytes
 - □ SRA = 5 petabytes
 - □ \$0.5-2 M/yr on S3



Democratized sequencing is here

- How to democratize expertise/HPC?
- What is the role of HPC in a peer-to-peer world?
- Lots of small groups doing compute
 - No experience with "proper" HPC usage
 - Jobs that run fine on a laptop don't translate to HPC



Acknowledgements



Karen Miga

genomeinformatics.github.io

- Adam Phillippy
- Arang Rhie
- Brian Walenz
- Alexander Dilthey



- ► ≒≒≒ Tim Smith
- M John Williams
- Sarah Kingan
- In Brittney Keel
- 🦙 Ben Rosen
- Sessional Petersen
- Mike Heaton
- Selection
 Selection
- Michael Hunkapiller
- Kerstin Howe
- Shane McCarthy
- Olivier Fedrigo
- UGP assembly group



