



A fully phased accurate assembly of an individual human genome

Tobias Marschall

Heinrich Heine University Düsseldorf, Germany
twitter: @tobiasmarschal

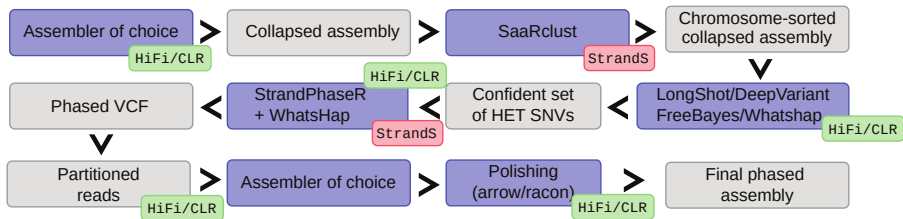
January 15, 2020
SMRT Informatics Developers Conference @ PAG

bioRxiv preprint first posted online Nov. 26, 2019; doi: <http://dx.doi.org/10.1101/855049>. The copyright holder for this preprint (which was not peer-reviewed) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

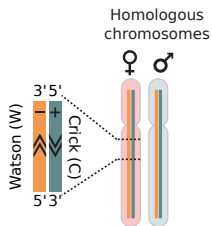
A fully phased accurate assembly of an individual human genome

David Porubsky^{1*}, Peter Ebert^{2*}, Peter A. Audano¹, Mitchell R. Vollger¹, William T. Harvey¹, Katherine M. Munson¹, Melanie Sorensen¹, Arvis Sulovari¹, Marina Haukness³, Maryam Ghareghani^{2,4}, Human Genome Structural Variation Consortium⁵, Peter M. Lansdorp^{6,7}, Benedict Paten³, Scott E. Devine⁸, Ashley D. Sanders⁹, Charles Lee¹⁰, Mark J.P. Chaisson¹¹, Jan O. Korbel⁹, Evan E. Eichler^{1,12†}, Tobias Marschall^{2,4†}

Workflow Overview

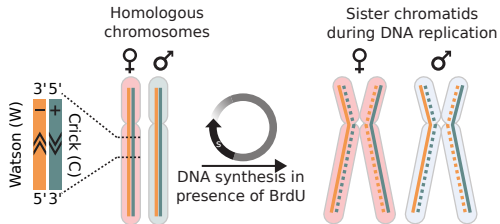


Strand-specific single-cell sequencing (Strand-seq)



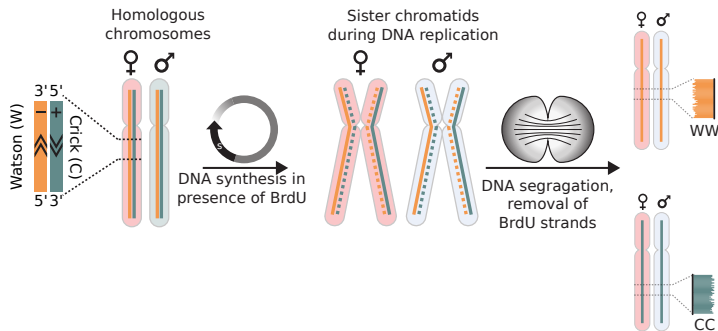
[Figure adapted from D Porubský, A Sanders, et al., Genome Research, 2016]

Strand-specific single-cell sequencing (Strand-seq)



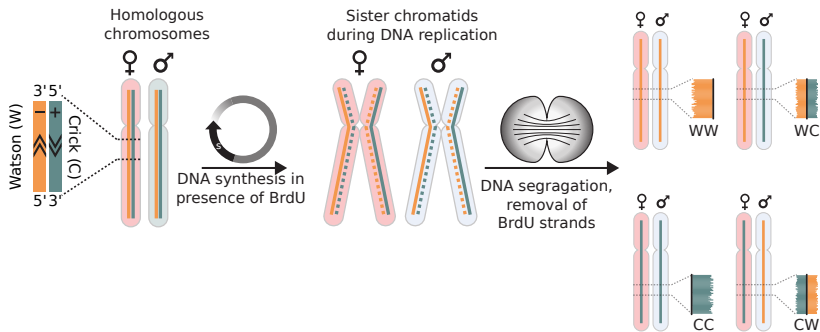
[Figure adapted from D Porubský, A Sanders, et al., Genome Research, 2016]

Strand-specific single-cell sequencing (Strand-seq)



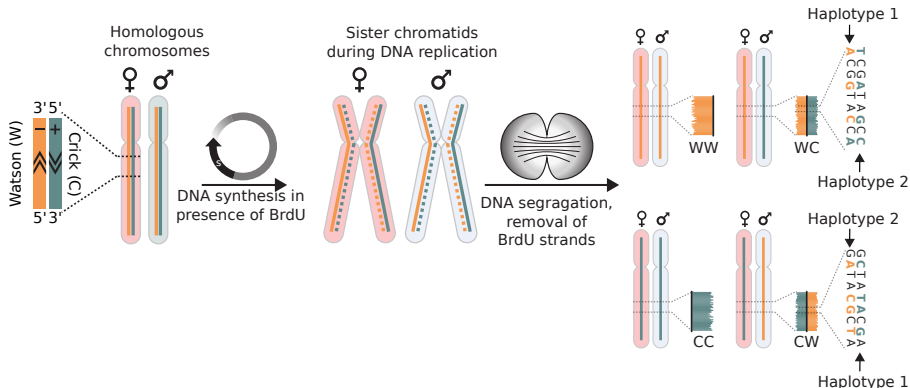
[Figure adapted from D Porubský, A Sanders, et al., Genome Research, 2016]

Strand-specific single-cell sequencing (Strand-seq)



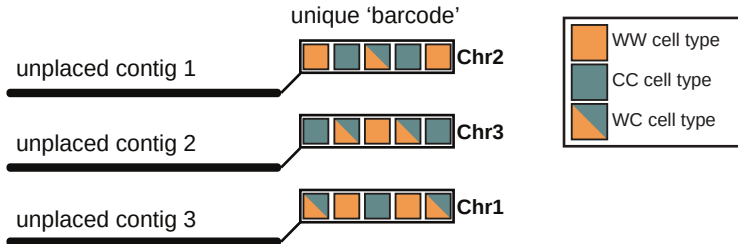
[Figure adapted from D Porubský, A Sanders, et al., Genome Research, 2016]

Strand-specific single-cell sequencing (Strand-seq)



[Figure adapted from D Porubský, A Sanders, et al., Genome Research, 2016]

Strand-seq can cluster contigs by chromosome



Bioinformatics, 34, 2018, i115–i123
doi: 10.1093/bioinformatics/bty290
ISMB 2018

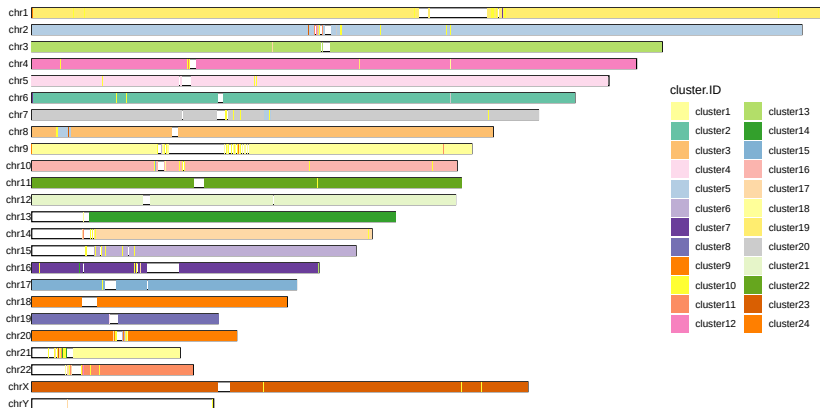
OXFORD

Strand-seq enables reliable separation of long reads by chromosome via expectation maximization

Maryam Ghareghani^{1,2,3,†}, David Porubský^{1,2,†}, Ashley D. Sanders⁴,
Sascha Meiers⁴, Evan E. Eichler^{5,6}, Jan O. Korbel⁴ and Tobias Marshall^{1,2,*}

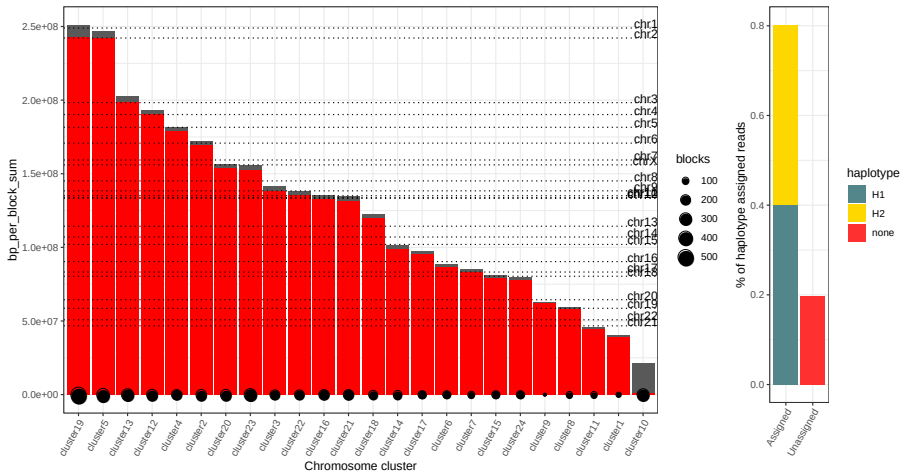
Clustering “collapsed” Canu assembly

HiFi reads of HG00733

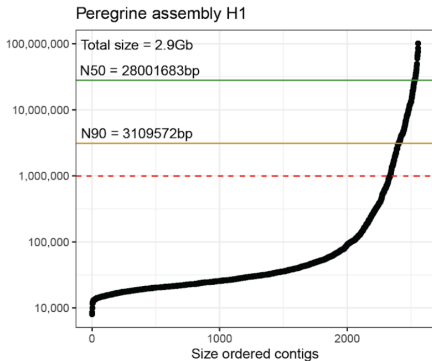


99.6% of contig sequence maps back to correct chromosome

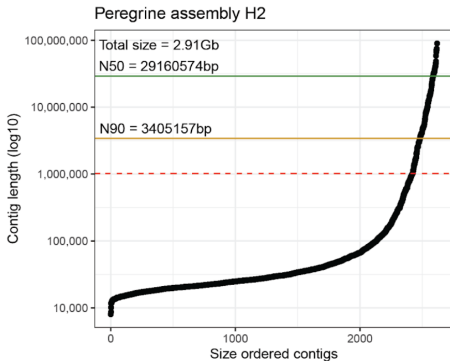
Phasing of each cluster



Resulting phased assembly



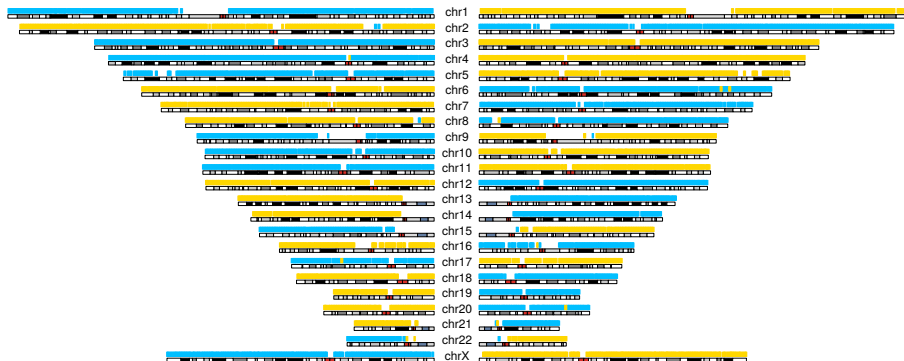
N50 after fixing assembly errors:
25,893,285 bp



N50 after fixing assembly errors:
28,927,707 bp

QV 40.8 as estimated from 77 BACs

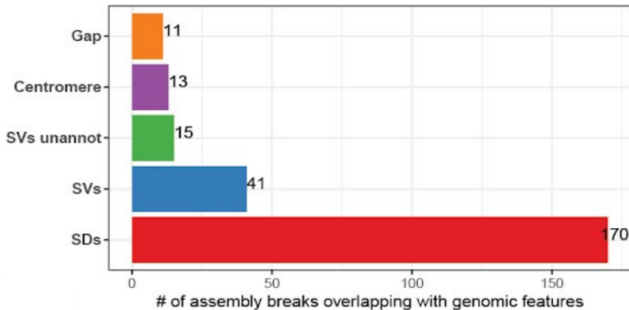
Phasing validation using parents



Switch error rate 0.4%

Universal Assembly Breaks

Assembly breaks shared by Peregrine and Canu



- 150/250 shared with trio-binned ONT assembly using Shasta
- 64/250 shared with T2T assembly of CHM13

Summary

- Partitioning reads without reference and parent information
- Haplotype-resolved assemblies N50 >25Mbp, QV >40
- Also works with CLR (not shown here)
- Universal assembly breaks shared across sequencing platforms and assembly tools

Funding:



National Institutes of Health



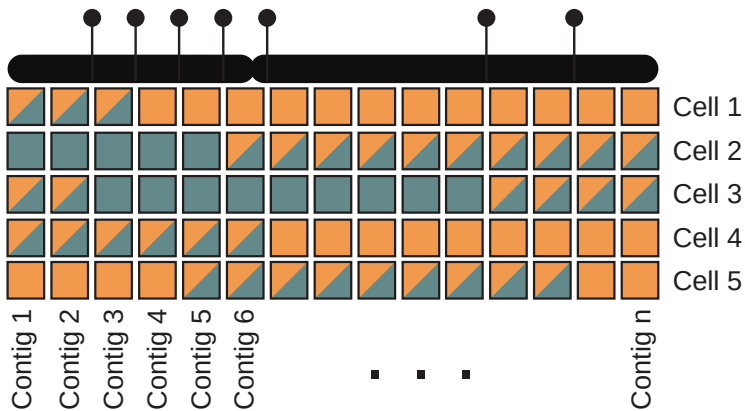
Deutsche
Forschungsgemeinschaft



Bundesministerium
für Bildung
und Forschung

Appendix

Strand-seq can cluster contigs by chromosome



Using Strand-seq to detect assembly errors

ctg	start	end	direction	cluster.ID
000008F	1	10772584	dir	cluster2
000008F	10772585	60481328	dir	cluster1
000006F	1	25346285	dir	cluster15
000006F	25346286	66162004	revcomp	cluster15

Chimerism



Misorientation

