Nanopore Sequencing

Sam Kovaka (Many slides by Michael Schatz)

Mar 2, 2020 Lecture 11: Applied Comparative Genomics



















The advantages of SMRT sequencing Roberts, RJ, Carneiro, MO, Schatz, MC (2013) *Genome Biology.* 14:405

Oxford Nanopore Technologies (ONT)



Nanopore Sequencing

Sequences DNA/RNA by measuring changes in ionic current as nucleotide strand passes through a pore



nanoporetech.com/applications/dna-nanopore-sequencing

Oxford Nanopore MinION





Nanopore Read Lengths

A typical MinION run produces ~10Gbp worth of reads of with a mean read length of ~10Kbp

 "Ultra-long" runs produce many reads > 100Kbp, with the longest read ever observed ~2.3Mbp



Nanopore sequencing and assembly of a human genome with ultra-long reads Jain et al. (2018) Nature Biotechnology: https://www.nature.com/articles/nbt.4060

Nanopore Read Quality

ONT reads typically have a mean error rate of ~10%

Predominantly indels



TTGTAAGCAGTTGAAAACTATGTGT <mark>G</mark> GATTTAG <mark>A</mark> ATAAAGAACATGAAAG
ATTATAAA-CAGTTGATCCATT-AGAAGA-AAACGCAAAAGGCGGCTAGG
CAACCTTGAATGTAATCGCACTTGAAGAACAAGATTTTATTCCGCGCCCCG
T <mark>A</mark> ACGAATC <mark>A</mark> AGATTCTGAAAACA <mark>C</mark> AT-AT <mark>AACA</mark> ACCTCCAAAA-CACAA
-AGGAGG <mark>GGA</mark> AA <mark>GGGGG</mark> GAATATCT-AT <mark>A</mark> AAAGATTACAAATT <mark>A</mark> GA-TGA
ACT-AATTCACAA <mark>T</mark> A-AATAACACTTTTA-ACA <mark>G</mark> AATTGAT-GGAA-GTT
TC <mark>G</mark> GAGAGATCC <mark>A</mark> AAACAAT <mark>G</mark> GGC-ATCG <mark>C</mark> CTTTGA-GTTAC-AATCAAA
ATCCAGT <mark>G</mark> GAAAATATA <mark>AT</mark> TTATGC <mark>A</mark> ATCCA <mark>G</mark> GAACTTATTCACAATTAG

Single Molecule Sequences



"Corrective Lens" for Sequencing



"Corrective Lens" for Sequencing



Consensus Accuracy and Coverage



Coverage can overcome random errors

- Dashed: error model from binomial sampling; Solid: observed accuracy
- Unfortunately, ONT has some non-random errors (mainly homopolymers)

$$CNSError = \sum_{i=\lceil c/2\rceil}^{c} \binom{c}{i} (1-e)^{n-i}$$

ONT Assembly Accuracy

Assemblies are quite contiguous, but percent identity maxes out ~99%

- Depends on organism basecallers are mainly trained on human
- Polishing from other technologies necessary for reference-quality



Benchmarking of long-read assemblers for prokaryote whole genome sequencing Wick and Holt (2019) F1000 Research: https://doi.org/10.12688/f1000research.21782.1



Translation of raw signal into basepairs



Translation of raw signal into basepairs

Early basecallers began by estimating k-mer boundaries using "events", which were then input to an HMM

Modern basecalers use neural networks directly on raw signal



(Based on probability of event matches)

ONT releases k-mer models with expected current distribution of every k-mer



DNA Base-Calling from a Nanopore Using a Viterbi Algorithm Timp et al. (2012) *Biophysical Journal*



Certain k-mers can be eliminated based on possible transitions



DNA Base-Calling from a Nanopore Using a Viterbi Algorithm Timp et al. (2012) *Biophysical Journal*



Final sequence determined by most probable k-mers



"DNA Base-Calling from a Nanopore Using a Viterbi Algorithm" Timp et al. (2012) *Biophysical Journal*

Basecaller/Pore Timeline

Development of both pore chemistry and basecalling algorithms is responsible for improvement in accuracy



From squiggle to basepair: computational approaches for improving nanopore sequencing read accuracy Rang *et al* (2018) *Genome Biology.* https://doi.org/10.1186/s13059-018-1462-9

Basecaller/Pore Timeline

Development of both pore chemistry and basecalling algorithms is responsible for improvement in accuracy



From squiggle to basepair: computational approaches for improving nanopore sequencing read accuracy Rang *et al* (2018) *Genome Biology.* https://doi.org/10.1186/s13059-018-1462-9

New Pore Chemistries

ONT recently released "R10" pore chemisity

- Two bottlenecks where nucleotides affect current
- Spans longer homopolymers
- May still have non-random errors, but profile is different



From 2018 London Calling Keynote

https://vimeo.com/272526835

More Throughput





MinION Quick Mobile Sequencing \$1k / instrument 5-8 GB / day

PromethION

High Throughput Desktop Sequencer \$75k / instrument >>1000GB / day

Nanopore Performance at CSHL

Sara Goodwin



Part of collaboration between JHU and CSHL to sequence 100 tomato genomes in 100 days

Nanopore Performance at CSHL

Sara Goodwin



Telomere-to-Telomere (T2T)

T2T consortium aims to finish the human genome

- Many gaps still exist, particularly around centromeres
- Uses ultra-long ONT reads, in addition to PacBio HiFi and other tech
- Finished chrX, close to finishing chr8, only 22 more to go



Karen Miga, London Calling 2019

Telomere-to-telomere assembly of a complete human X chromosome Miga et al. (2019) *BioRxiv.* https://doi.org/10.1101/735928

DNA Modification Detection

ONT can detect methylation from raw signal

• Or any other modification that changes ionic current



Piercing the dark matter: bioinformatics of long-range sequencing and mapping Sedlazeck et al. (2018) *Nature Reviews Genetics.* 19:329

Nanopore Direct RNA-seq



ONT can sequence RNA molecules directly

cDNA sequencing erases modifications and structure

Direct RNA-seq has potential to read both



Long-Read RNA-seq Assembly

Both direct and cDNA ONT RNA-seq produce some transcript fragments

• I helped develop StringTie2, which applies short-read transcriptome assembly methods to long-reads





Transcriptome assembly from long-read RNA-seq alignments with StringTie2

Sam Kovaka, Aleksey V. Zimin, Geo M. Pertea, Roham Razaghi, Steven L. Salzberg & Mihaela Pertea <u>Genome Biology</u> 20, Article number: 278 (2019) Cite this article

VolTRAX - Library Prep (+ sequencing?)



- 3.6 kb DNA, standard ligation library preparation
- · Sample + pores in one droplet
- · Pores inserted, then library sequenced
- Droplet size 10nL, could be 4.5 nL with current chip







Proof of concept array demonstrated

- No crosstalk data taken directly from cartridge
- Wide range of experiments possible
- Will include MinKNOW control and feedback
- Data being collected for model training.



Extremely Portable Sequencing!



Kate Rubins sequencing DNA on the ISS

Ebola Surveillance

LETTER

doi:10.1038/nature16996

Real-time, portable genome sequencing for Ebola surveillance

Joshua Quick1*, Nicholas J. Loman1*, Sophie Duraffour2.1*, Jared T. Simpson4.1*, Ettore Severif*, Lauren Cowley7*, Joseph Akoi Bore², Raymond Koundouno², Gytis Dudas⁸, Amy Mikhail⁷, Nobila Ouédraogo⁹, Babak Afrough^{2,10} Amadou Bah2,11, Jonathan H. J. Baum2,3, Beate Becker-Ziaja2,5, Jan Peter Boettcher2,12, Mar Cabeza-Cabrerizo2,3, Álvaro Camino-Sánchez⁷, Lisa L. Carter^{2,13}, Juliane Doerrbecker^{2,3}, Theresa Enkirch^{2,14}, Isabel García-Dorival^{2,13}, Nicole Hetzelt^{2,12}, Julia Hinzmann^{2,22}, Tobias Holm^{2,3}, Liana Eleni Kafetzopoulou^{3,19}, Michel Koropogul^{2,37}, Abigael Kosgey^{2,18}, Eeva Kuisma^{2,20}, Christopher H. Logue^{2,20}, Antonio Mazzarelli^{2,19}, Sarah Meisel^{2,3}, Marc Mertens^{2,20}, Janine Michel^{2,12}, Didier Ngabo^{2,35}, Katja Nitzsche^{2,3}, Elisa Pallasch^{2,3}, Livia Victoria Patrono^{2,3}, Jasmine Portmann^{2,21}, Johanna Gabriella Repits^{2,22}, Natasha Y. Rickett^{2,15,23}, Andreas Sachse^{2,12}, Katrin Singethan^{2,24}, Inés Vitoriano^{2,10}, Rahel L. Yemanaberhan^{2,3}, Elsa G. Zekeng^{2,15,23}, Trina Racine²⁵, Alexander Bello²⁵, Amadou Alpha Sall²⁶, Ousmane Faye²⁶, Oumar Faye²⁶, N'Faly Magassouba²⁷, Cecelia V. Williams^{28,29}, Victoria Amburgey^{28,29}, Linda Winona^{28,29}, Emily Davis^{29,30}, Jon Gerlach^{29,30}, Frank Washington^{29,30}, Vanessa Monteil³¹, Marine Jourdain³¹, Marion Bererd³¹, Alimou Camara³¹, Hermann Somlare³¹, Abdoulaye Camara¹⁰, Marianne Gerard¹¹, Guillaume Bado¹¹, Bernard Baillet¹², Déborah Delaune^{12,23}, Koumpingnin Yacouba Nebie³⁴, Abdoulaye Diarra³⁴, Yacouba Savane³⁴, Raymond Bernard Pallawo³⁴, Giovanna Jaramillo Gutierrez³⁵, Natacha Milhano^{6,36}, Isabelle Roger³⁴, Christopher J. Williams^{6,07}, Facinet Yattara¹⁷, Kuiama Lewandowski⁷⁰, James Taylor³⁸, Phillip Rachwal³⁸, Daniel J. Turner³⁹, Georgios Pollakis^{15,23}, Julian A. Hiscox^{15,23}, David A. Matthews⁴⁰, Matthew K. O'Shea⁴⁵, Andrew McD. Johnston²⁶, Duncan Wilson⁴⁷, Emma Hutley⁴², Erasmus Smilt⁴³, Antonino Di Caro^{2,19}, Roman Wölfel^{2,44}, Kilian Stoecker^{2,44}, Erna Fleischmann^{2,44}, Martin Gabriel^{2,3}, Simon A. Weller¹⁸, Lamine Kolvogui⁴⁵, Boubacar Diallo¹⁴, Sakoba Keita¹⁷, Andrew Rambaut^{8,46,47}, Pierre Formenty³⁴, Stephan Günther^{2,3} & Miles W. Carroll^{2,30,48,49}

Ebola Surveillance

LETTER

doi:10.1038/nature16996

Real-time, portable genome sequencing for Ebola surveillance

Joshua Quick1*, Nicholas J. Loman1*, Sophie Duraffour2.1*, Jared T. Simpson1.1*, Ettore Sev Joseph Akoi Bore², Raymond Koundouno², Gytis Dudas⁸, Amy Mikhail⁷, Nobila Ouedraog Amadou Bah2,11, Jonathan H. J. Baum2,3, Beate Becker-Ziaia2,5, Jan Peter Boettcher2,12, Mar Álvaro Camino-Sánchez⁷, Lisa L. Carter^{2,13}, Juliane Doerrbecker^{2,3}, Theresa Enkirch^{2,14}, Is Nicole Hetzelt^{2,12}, Julia Hinzmann^{2,12}, Tobias Holm^{2,3}, Liana Eleni Kafetzopoulou^{2,18}, Miche Eeva Kuisma230, Christopher H. Logue230, Antonio Mazzarelli239, Sarah Meisel23, Marc M Didier Ngabo^{2,35}, Katja Nitzsche^{2,3}, Elisa Pallasch^{2,3}, Livia Victoria Patrono^{2,3}, Jasmine Port Natasha Y. Rickett^{2,15,23}, Andreas Sachse^{2,12}, Katrin Singethan^{2,24}, Inés Vitoriano^{2,10}, Rahel Elsa G. Zekeng^{3,15,23}, Trina Racine²³, Alexander Bello²⁵, Amadou Alpha Sall²⁶, Ousmane Fa N'Faly Magassouba27, Cecelia V. Williams28,29, Victoria Amburgey28,29, Linda Winona28,29, Er Frank Washington^{29,30}, Vanessa Monteil³¹, Marine Jourdain³¹, Marion Bererd³¹, Alimou Cam Abdoulaye Camara²¹, Marianne Gerard²¹, Guillaume Bado²¹, Bernard Baillet²², Déborah Dela Abdoulaye Diarra34, Yacouba Savane34, Raymond Bernard Pallawo34, Giovanna Jaramillo Gu Isabelle Roger³⁴, Christopher J. Williams^{6,37}, Facinet Yattara¹⁷, Kuiama Lewandowski¹⁰, Jame Daniel J. Turner³⁰, Georgios Pollakis^{13,23}, Julian A. Hiscox^{13,23}, David A. Matthews⁴⁰, Matthew Andrew McD. Johnston⁴², Duncan Wilson⁴², Emma Hutley⁴², Erasmus Smit⁴³, Antonino Di G Kilian Stoecker^{2,44}, Erna Fleischmann^{2,44}, Martin Gabriel^{2,3}, Simon A. Weller¹⁸, Lamine Kolv Sakoba Keita¹⁷, Andrew Rambaut^{8,46,47}, Pierre Formenty³⁴, Stephan Günther^{2,3} & Miles W.O.



Figure 1 | Deployment of the portable genome surveillance system in Guinea. a, We were able to pack all instruments, reagents and disposable consumables within aircraft baggage. b, We initially established the genomic surveillance laboratory in Donka Hospital, Conakry, Guinea. c, Later we moved the laboratory to a dedicated sequencing laboratory in Coyah prefecture. d, Within this laboratory we separated the sequencing instruments (on the left) from the PCR bench (to the right). An uninterruptable power supply can be seen in the middle that provides power to the thermocycler. (Photographs taken by J.Q. and S.D.)

Ebola Surveillance



Figure 2 | Real-time genomics surveillance in context of the Guinea Ebola virus disease epidemic. a, Here we show the number of reported cases of Ebola virus disease in Gainea (red) in relation to the number of EBOV new patient samples (n = 137, in blae) generated during this study. b, For each of the 142 sequenced samples, we show the relationship between sample collection date (red) and the date of sequencing (blae). Twenty-eight samples were sequenced within three days of the sample being taken, and sixty-eight samples within a week. Larger gaps represent retrospective sequencing of cases to provide additional epidemiological context.







related to cases identified in Serra Leone. Samples are frequently clustered by grography Coldicated by colour of circle's and this provides information as to origins of new introductions, such as in the built epidemic in May 2015. Map figure adapted from KimpleMaps website (http://simplemaps. sint/recordering.gc).

COVID-19 () Surveillance

200 Oxford Nanopore sequencers have left UK for China, to support rapid, near-sample coronavirus sequencing for outbreak surveillance





ARTIC Network @NetworkArtic - Feb 5

First #nCoV2019 genome sequenced using the ARTIC @nanopore protocol and primer scheme by @Scalene is up on #GISAID. Sequenced by the Hangzhou CDC with assistance from @nanopore.artic.network/ncov-2019

Requesting technology Assessing method	Nampon Medity articlescentry and Service										
Coverage:	x 20005a (avendari)										
Fulling of second second											
Originating lab:	Heights, Carks for Dease	Control and Prevant	100								
Address: St Wogen Rust, designer flamtd 20000 Sample Bigheen by Tex sample presider Sample St gleen by Tex sample Samples Control for Disease Control and Presenter. Address: St Strapping Control 20000 Samples Control for Disease Control 20000 Samples Control for Disease Control 20000 Samples Control 20000											
						Alampia ID groom by the automitting laboratory	Beld's Margins Addition (1998) Janus, massa Rang, marina singling Max Kolan Ita Jina Ran, Brasil Rang, Kin Gan, Bharlang Chen, Kadin Rang				
						Adure					
						0.3	13 42	C	105	1	





james hadfield @hamesiadfield

When we started developing RAMPART the idea that you could analyse seq data in real time seemed rather aspirational. Today: a still-sequencing **#SARSCoV2** @nanopore genome from Brazil already up on @nextstrain. Amazing to get to work with @jaquegj @CaddeProject @NetworkArtic et al



https://nanoporetech.com/about-us/news/novel-coronavirus-ncov-2019covid-19-information-and-updates
Less Throughput



Flongle

- An adapter for MinION for smaller tests or experiments
- · Single-use, on-demand, cost efficient sequencing
- Suitable for quality checks, amplicons, smaller genomes, targeted regions, or those interested in diagnostics/other tests
- MinIT available to support IT/software needs



SmidgION (coming "soon")

- Designed to be our smallest sequencing device so far
- Same nanopore sensing technology as MinION and PromethION
- Designed for use with a smartphone in any location

https://nanoporetech.com/products

Targeted Sequencing

Often you're only interested in certain sequences

- Can be challenging to reach sufficient coverage with low yield
- For example: **pathogen DNA** enrichment or **targeting genes**



PCR doesn't work work well for ONT sequencing

• Limits read length and erases epigenetc modifications

CIRSPR/Cas9 Enrichment

Uses Cas9 to cut DNA around target region, then binds adapters to phosphorylated ends



Targeted Nanopore Sequencing with Cas9 for studies of methylation, structural variants, and mutations

Timothy Gilpatrick, Isac Lee, James E. Graham, Etienne Raimondeau, Rebecca Bowen, Andrew Heron, ⁽¹⁾ Fritz J Sedlazeck, ⁽¹⁾ Winston Timp

ReadUntil Sequencing

ONT devices can selectively eject reads in real-time



Nanopore Sequencing



Nanopore Sequencing with ReadUntil

Enables targeted sequencing without addition sample prep

• Requires rapid real-time read identification

MinION has up to 512 active channels, each reading 450 bp/sec

• Actual number of active channels is variable

ReadUntil Sequencing Reference On-Target Off-target Emichment Target Depletion **Off-target Reference Off-Target**

Utility for Nanopore Current ALignment to Large Expanses of DMAINCALLED

A novel streaming algorithm which maps raw nanopore signal as it is being sequenced

- Can map reads from all active MinION channels to reference tens of megabases in size
- Uses the mapping results to make ReadUntil decisions



Reference Genome Raw Signal TGCAAGCATGC1 FM Index Events 2 Log K-Mer Match Probabilities 5 -1.8 AT -2.4 TT -3.1 GC 6 -2.2 CA -1.6 TG -2.6 TG 7 -3.4 CT -3.7 AT 8 TC -2.8 CT -1.6 TC -4.6 AA 9 Alignment Forest 10 11 12 bioRχ SERVER FOR BIOLOG

Targeted nanopore sequencing by real-time mapping of raw electrical signal with UNCALLED

Sam Kovaka, Yunfan Fan, Bohan Ni, 💿 Winston Timp, 💿 Michael C Schatz doi: https://doi.org/10.1101/2020.02.03.931923

UNCALLED Signal Processing

- Stretches of similar signal are collapsed into events
 - Averages out noise and reduces amount of signal to process
 - Ideally each event represents a single k-mer, but many errors occur (50% stays, 1% skips)
- Probability of events matching each k-mer is then computed
 - Expected current for each k-mer modeled by normal distribution
 - ONT releases 6-mer models (I use 5-mers)



FM Index

- Used by many aligners such as BWA, Bowtie, and HISAT
- UNCALLED uses BWA's FM index
 - Interchangeable started with my own implementation



FM Index Search



Size of range = number of possible alignments

FM Index Search w/ Ambiguity



FM Index Search w/ Ambiguity



UNCALLED Algorithm



Mapping Reads to E. coli

Mapped 100K E. coli reads to the E. coli reference genome

- Running on a single core of a 3.0 GHz Intel Xeon Gold 6136
- Estimated accuracy using minimap2 as ground truth



Mapping Reads to E. coli

Mapped 100K E. coli reads to the E. coli reference genome

- Running on a single core of a 3.0 GHz Intel Xeon Gold 6136
- Estimated accuracy using minimap2 as ground truth



	Ρ	Ν
Т	85.70%	7.99%
F	1.03%	5.28%

Mapping Reads to E. coli

Mapped 100K E. coli reads to the E. coli reference genome

- Running on a single core of a 3.0 GHz Intel Xeon Gold 6136
- Estimated accuracy using minimap2 as ground truth



75.4% of "FPs" were not aligned by Minimap2

92.4% of FPs that Minimap2 aligned are explained by repeats

UNCALLED on Zymo

Tested UNCALLED on the Zymo high molecular weight mock microbial community standard

- Contains seven bacteria and one fungus
- Mapped reads to the full 41Mbp reference at a rate of ~6,000 bp/thread/sec
- Match minimap2 alignments with 96% accuracy

Sequencing by Yunfan Fan



Read lengths:

Median: 12.2Kbp Mean: 15.9Kbp N50: 24.7Kbp

- Mapping to a reference of the all bacteria
- Ejecting reads if they map within the first 4,500bp (10 sec), enriching fungal sequence
- Sequenced same sample on flowcell of similar quality directly next to ReadUntil run as control
- Running with 48 threads on a 3.0 GHz Intel Xeon Gold 6136 (probably overkill)

- Mapping to a reference of the all bacteria
- Ejecting reads if they map within the first 4,500bp (10 sec), enriching fungal sequence
- Sequenced same sample on flowcell of similar quality directly next to ReadUntil run as control
- Running with 48 threads on a 3.0 GHz Intel Xeon Gold 6136 (probably overkill)

Read decision accuracy:

Correctly eject: 3	, 467 , 492	96.74%
Correctly kept:	49,282	99.78%

- Mapping to a reference of the all bacteria
- Ejecting reads if they map within the first 4,500bp (10 sec), enriching fungal sequence
- Sequenced same sample on flowcell of similar quality directly next to ReadUntil run as control
- Running with 48 threads on a 3.0 GHz Intel Xeon Gold 6136 (probably overkill)

Read decision accuracy:

Correctly eject: 3	, 467, 492	96.74%
Correctly kept:	49,282	99.78%



- Mapping to a reference of the all bacteria •
- Ejecting reads if they map within the first • 4,500bp (10 sec), enriching fungal sequence
- Sequenced same sample on flowcell of • similar quality directly next to ReadUntil run as control
- Running with 48 threads on a 3.0 GHz • Intel Xeon Gold 6136 (probably overkill)

Read decision accuracy:

Correctly eject: 3,	467,492	96.74%
Correctly kept:	49,282	99.78%



Fold change:

4.50

Human Gene Enrichment

Enriching GM12878 for genes associated with hereditary cancer

- Could provide a low-cost and portable method to identify, SNVs, SVs, and epigenetic modifications
- Mapping to a 18.6 Mbp reference, ejecting reads that don't map in 3 sec, ran for 72 hours
- Also masked low-complexity and repetitive regions, and performed two nuclease flushes on each run



All 148 genes from every Invitae cancer panel:

	•		•			-
ABRAXAS1 ATP	CDKN1B CDKN1C	EXT1 EZH2	HOXB13 HRAS	NOP10 NTHL1	RECQL RECOL4	SDHC XRCC2 SDHB
AKT1	CEBPA	FANCA	KIF1B	PALB2	REST	SDHD
ALK	CDKN2A	FANCB	KIT	PALLD	RET	SLX4
AP2S1	CEP57	FANCC	LZTR1	PARN	RINT1	SMAD4
APC	CHEK2	FANCD2	MC1R	PDGFRA	RNF43	SMARCB1
ATM	CFTR	FANCE	MAX	PH0X2B	RPL11	SMARCA4
ATR	CPA1	FANCF	MEN1	PIK3CA	RPL15	SPINK1
AXIN2	CTNNA1	FANCG	MET	PMS2	RPL26	SMARCE1
BAP1	CTC1	FANCI	MITF	POLD1	RPL35A	STK11
BARD1	DIS3L2	FANCL	MLH1	P0T1	RPL5	SUFU
BMPR1A	CTRC	FANCM	MLH3	POLE	RPS10	TERT
BLM	CTR9	FLCN	MRE11	PRKAR1A	RPS19	TERC
BRCA1	DICER1	FH	MSH2	PRSS1	RPS20	TINF2
BRCA2	DKC1	GALNT12	MSH3	PTCH1	RPS24	TMEM127
CASR	EGLN1	GATA1	MSH6	PTCH2	RPS26	TP53
BUB1B	EGFR	GEN1	NBN	PTEN	RPS7	TSC1
BRIP1	EPCAM	GATA2	MUTYH	RAD50	RTEL1	TSC2
CDC73	ENG	GNA11	NF1	RAD51C	RUNX1	VHL
CDK4	ERCC4	GPC3	NF2	RAD51D	SDHA	WRN
CDH1	EXT2	GREM1	NHP2	RB1	SDHAF2	WT1

Hereditary Cancer Gene Enrichment



Hereditary Cancer Gene Enrichment



Hereditary Cancer Gene Analysis

Small variant calling with Clair

	Dataset	Precision	Recall	F1
SNPs Indels	Control (5x)	41.9%	39.4%	0.406
SNPs	UNCALLED (29.6x)	92.8%	97.6%	0.951
	ONT WGS (51.1x)	93.2%	98.5%	0.958
	Control (5x)	37.6%	23.4%	0.288
Indels	UNCALLED (29.6x)	80.4%	73.1%	0.766
	ONT WGS (51.1x)	79.9%	72.7%	0.761

Hereditary Cancer Gene Analysis

Small variant calling with Clair

	Dataset	Precision	Recall	F1
	Control (5x)	41.9%	39.4%	0.406
SNPs	UNCALLED (29.6x)	92.8%	97.6%	0.951
	ONT WGS (51.1x)	93.2%	98.5%	0.958
	Control (5x)	37.6%	23.4%	0.288
Indels	UNCALLED (29.6x)	80.4%	73.1%	0.766
	ONT WGS (51.1x)	79.9%	72.7%	0.761

Structural variant calling with Sniffles

- 50 SVs detected from UNCALLED reads
- 100% concordance with 50x coverage ONT WGS and 30x PacBio HiFi reads
- More than double the number of SVs detected with 50x coverage whole-genome Illumina sequencing

Hereditary Cancer Gene Analysis

Small variant calling with Clair

	Dataset	Precision	Recall	F1
	Control (5x)	41.9%	39.4%	0.406
SNPs	UNCALLED (29.6x)	92.8%	97.6%	0.951
	ONT WGS (51.1x)	93.2%	98.5%	0.958
	Control (5x)	37.6%	23.4%	0.288
Indels	UNCALLED (29.6x)	80.4%	73.1%	0.766
	ONT WGS (51.1x)	79.9%	72.7%	0.761

Structural variant calling with Sniffles

- 50 SVs detected from UNCALLED reads
- 100% concordance with 50x coverage ONT WGS and 30x PacBio HiFi reads
- More than double the number of SVs detected with 50x coverage wholegenome Illumina sequencing

Methylation calling with Nanopolish

						2-1				w -1 w	
		100.00	-			and the second s			and a state	- anima	
ACALLED Cor	10.0	10.2.00	-								
	111	1.1						11		1 1 1 1	1 () ((W
NOALLED Reen			1			1 1 1 1 1 1 1					
INT WEB Cav											
	1. 1				-		inii)		ii n i		• 111 sm
n/1 mülti faadt	11	111				1.10.000					
agains -		Test International Property in which the local division of the loc	-	and the second second	AT AD IN THE OWNER	series and the second of the	a rest to be set the		-	and the second second	

.96 Pearson correlation coefficient with 50x ONT WGS sequencing

Found evidence of Xinactivation

Structural Variant Breakdown

56 concordant SVs overall

• 39 insertions, 17 deletions

Classified insertions and deletions by repeat alignment and overlap

- 56% of insertions and 41% of deletions are simple repeats or low-complexity
- 16% of insertions and 29% of deletions overlap or align to Alu elements

One SV was found in an exon

- Heterozygous Alu insertion in MUTYH
- Mutations in this gene promote colorectal and breast cancers

	, 10 1 ,	11 p32.2 p32.1 p32 q32 q32 q32 q32	981 911
	29,800 hp	669 bp	30,100 hp
UNCALLED SVs			
UNCALLED Cov	CONTRACTOR OF A	and the state of the	
		III	
UNCALLED Reads			
ONT WCS SVe			
ONT WGS SVS	And in case of the local division of the loc		Comment of the
ONT WGS Reads		IM	
		III	
PacBio SVs			
PacBio Cov			
PacBio Reads			
Illusian Che		ADA	
Illumina SVS		No variants Found	-
Illumina Reads	1		
Sequence 🗕			COLUMN TWO IS NOT
Genes		MUTYH	

Class Project

- UNCALLED started as a project for this class!
- How we split the work between the three of us:
 - Collecting/parsing raw nanopore signal data
 - Signal processing/k-mer matching
 - FM Index construction/basic search algorithm
- All of us brainstormed how the algorithm should work
- We did not have a functional aligner in the end
 - Created a signal-based FM index (later turned out to be unnecessary)
 - Figured out how to compute event/k-mer match probabilities (but messed up signal normalization)
 - Could produce seed alignments based on a very simple algorithm (but had no way to filter the many many false positives)
- Despite the incompleteness it was a successful project!



Questions?

Reduced UNCALLED Yield

UNCALLED has ~50% the yield of control

- Fewer pores are active throughout
- Partially explained by ejections causing more gaps between reads, but not entirely
- Pores do not seem to be "dying" faster

Likely due to ejections causing more pore blockages

• Maybe single stranded DNA self binding, or more reads mean more chances to clog



Nuclease Flush Results



Lower threshold means

• Takes more time to process the event



Lower threshold means

- Takes more time to process the event
- Higher chance of finding the right k-mer



Lower threshold means

- Takes more time to process the event
- Higher chance of finding the right k-mer



Lower threshold means

- Takes more time to process the event
- Higher chance of seeing correct k-mer



Smaller FM index ranges represent fewer genomic locations

- Fewer possible distinct sequences
- Less potential for branching

Threshold gets lower (considers more k-mers) as FM index range gets smaller

• Must change with reference size/repeat content

