



PACIFIC
BIOSCIENCES®

PacBio SMRT Technology & Long Read Sequencing Applications Overview

Sonny S. Mark, Ph.D.
Staff Scientist, Pacific Biosciences

April. 11. 2019

Agenda

1. SMRT Sequencing Technology & Sequel System Overview

2. PacBio Long-Read Sequencing Applications

a. Whole Genome Sequencing

- De Novo Assembly
- Structural Variation Detection

b. Targeted Sequencing

c. Analysis of Complex Populations (Metagenomics, Cancer)

d. (RNA) Transcript Isoform Sequencing (Iso-Seq Method)

e. Epigenetics

3. Summary

***A PacBio Local
SMRT Grant
Award Program at
McMaster!***



PACIFIC
BIOSCIENCES®



SMRT Sequencing Technology & Sequel System Overview

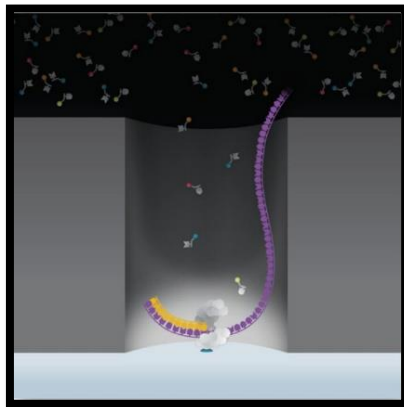
Sequence with Confidence. Advance genomics with long-read sequencing, enabled by single molecule real-time sequencing.

SINGLE MOLECULE, REAL-TIME SEQUENCING

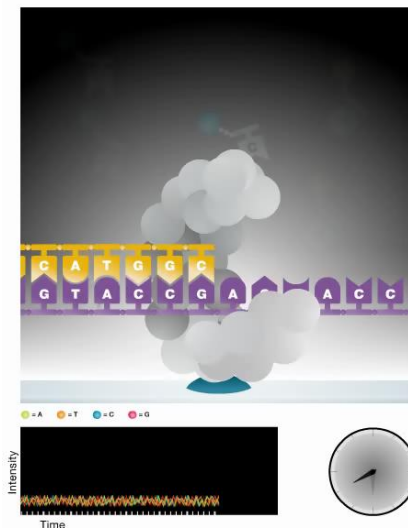
SMRT Sequencing Technological Innovations Enable Long Reads from Single Molecules

1. A single polymerase-bound DNA molecule is immobilized within a nanofabricated zero-mode waveguide (ZMW)
2. Harnessing the power of natural DNA synthesis, phospho-linked fluorescent-labeled nucleotide analogs are incorporated into the DNA molecule by the polymerase
3. The synthesis process is observed in real-time across 1,000,000 ZMWs on a single SMRT Cell consumable by converting the detected light signals into base calls

Zero-mode Waveguide (ZMW)



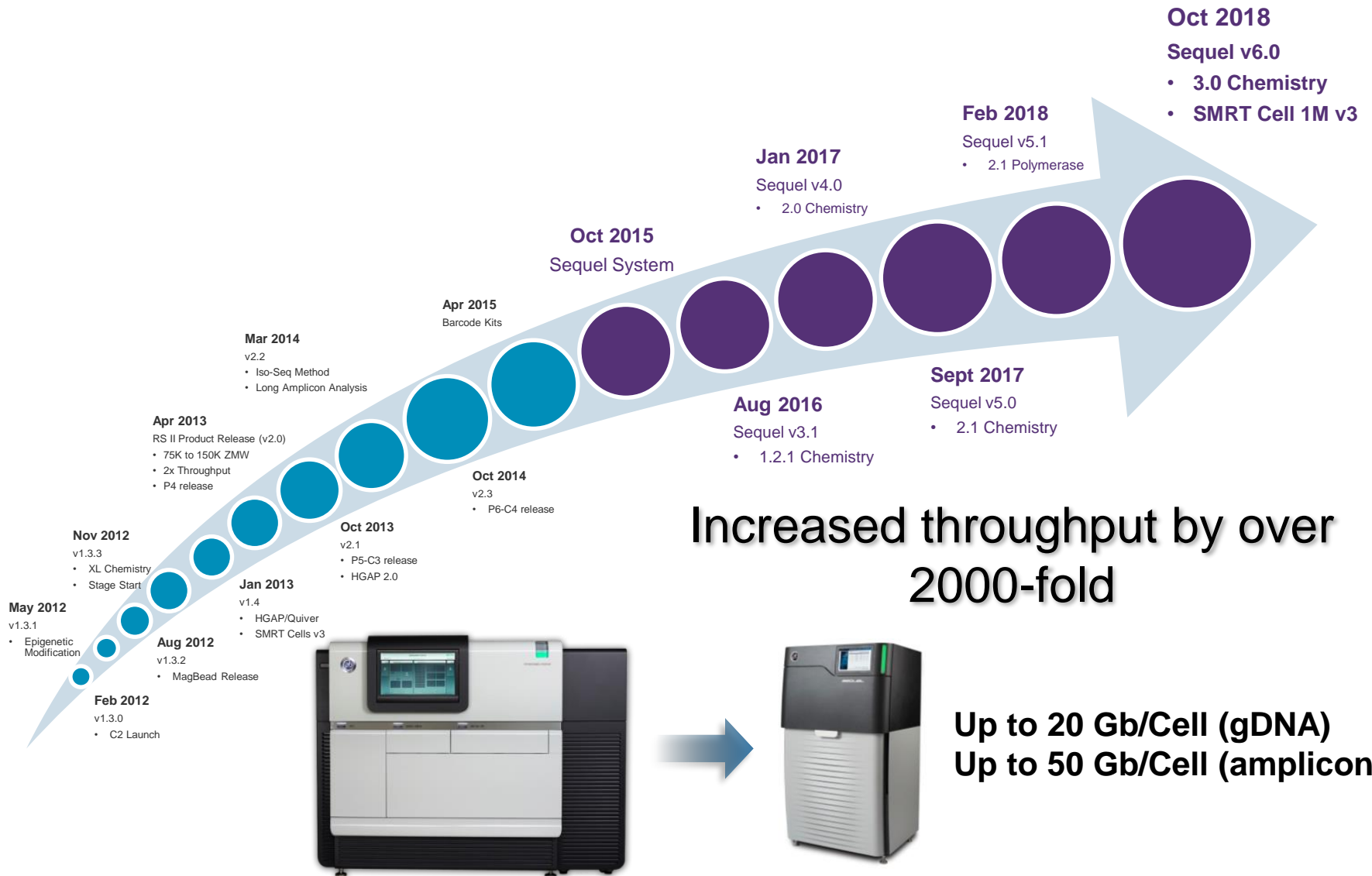
Phospho-linked Nucleotides



SMRT Cells with 10⁶ ZMWs

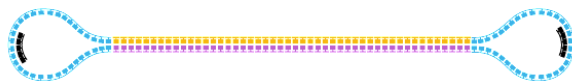


PACBIO PRODUCT RELEASES OVER THE LAST SEVEN YEARS



TWO APPROACHES FOR SMRT SEQUENCING

Standard Sequencing for Single-Pass Continuous Long Reads (CLR)

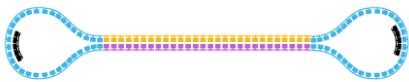


Large Insert Sizes
(Approx. >15 kb)



Generates **one pass** on each library insert molecule sequenced

Circular Consensus Sequencing (CCS)



Small Insert Sizes
(Approx. ≤15 kb)

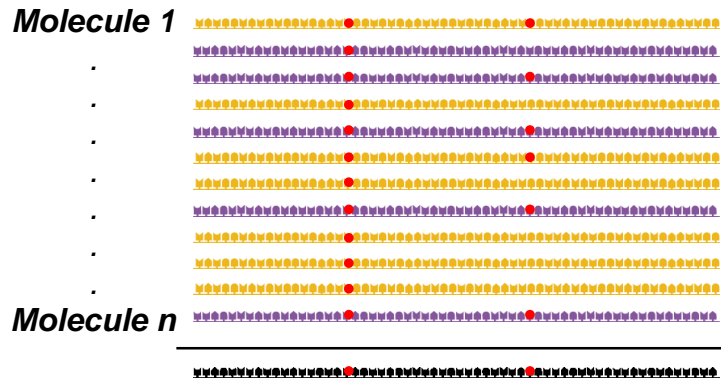
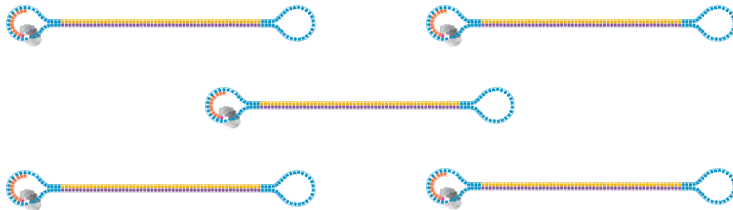


Continued generation of reads per insert size

Generates **multiple passes** on each library insert molecule sequenced

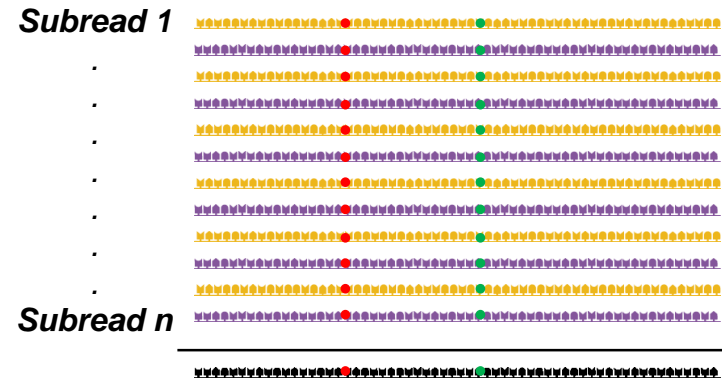
CONSENSUS SEQUENCE GENERATION FROM MULTIPLE INDIVIDUAL READS (STANDARD CLR SEQUENCING) OR FROM MULTIPLE SUBREADS OF THE SAME DNA MOLECULE (CCS)

Standard Sequencing (CLR)



Consensus sequence

Circular Consensus Sequencing (CCS)



Consensus sequence

Example Application:

- WGS *de novo* assembly
- Structural variation detection
- Base modification detection

Example Application:

- Minor variant detection (viral, cancer)
- Full-length 16S metagenomics
- Full-length transcriptomics (Iso-Seq analyses)

THE SMRT SEQUENCING ADVANTAGE: KEY PERFORMANCE CRITERIA

SMRT Sequencing can *simultaneously* provide:

Long Reads

- Average read lengths many tens of kb

High Accuracy

- Free of systematic errors
- Achieves >99.999% (Q50)

Uniform Coverage

- Least GC content and sequence complexity bias

Single-Molecule Resolution

- Long reads with high *single-molecule* accuracy
- Resolve complex mixtures

Epigenetic Detection

- Characterize epigenome
- No separate sample preparation required

SMRT Sequencing Advantages



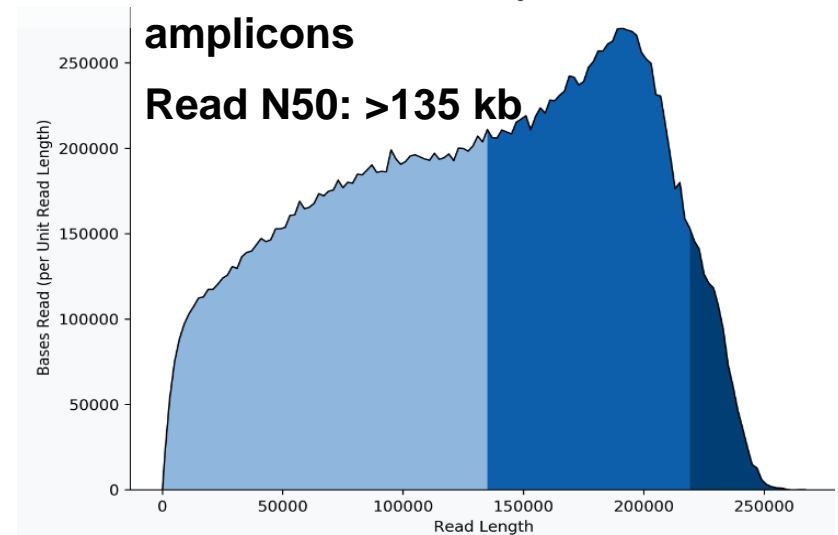
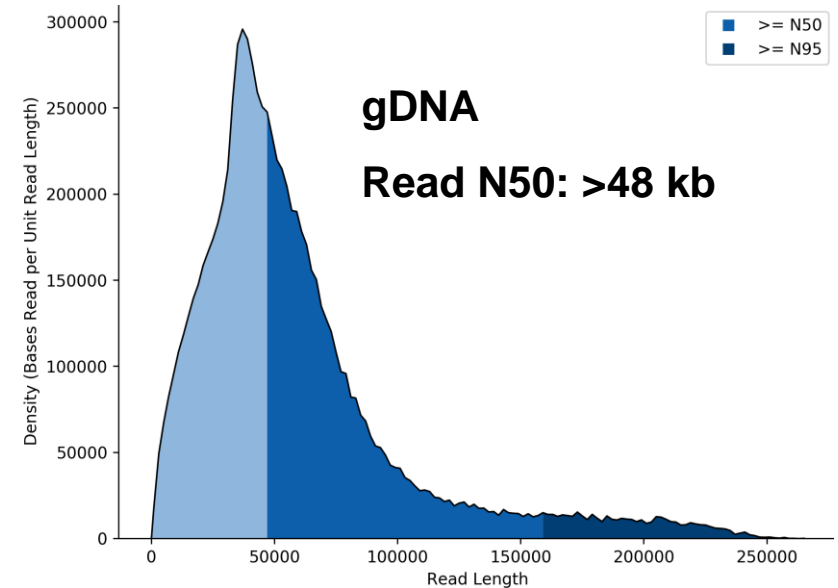
SEQUENCING PERFORMANCE (SEQUEL SYSTEM 3.0 CHEMISTRY)

Long Reads

- Up to 30 kb average (gDNA)
- Up to 100 kb average (amplicons)

High Consensus Accuracy

Uniform, Unbiased Coverage



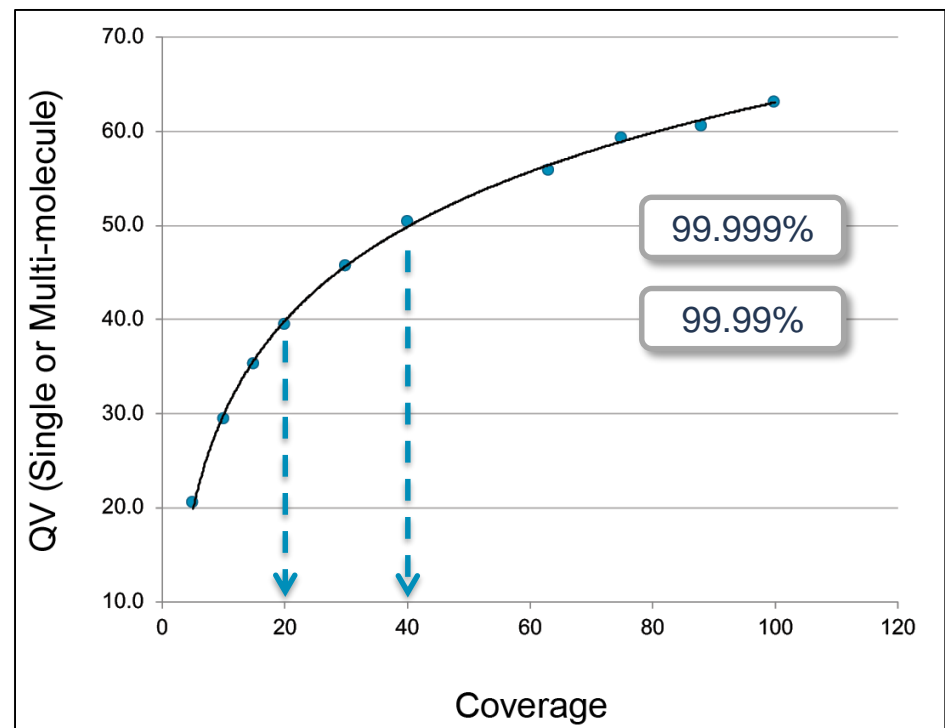
SEQUENCING PERFORMANCE (SEQUEL SYSTEM 3.0 CHEMISTRY)

Long Reads

High Consensus Accuracy

- >QV40 at 20-fold
- >QV50 at 40-fold

Uniform, Unbiased Coverage



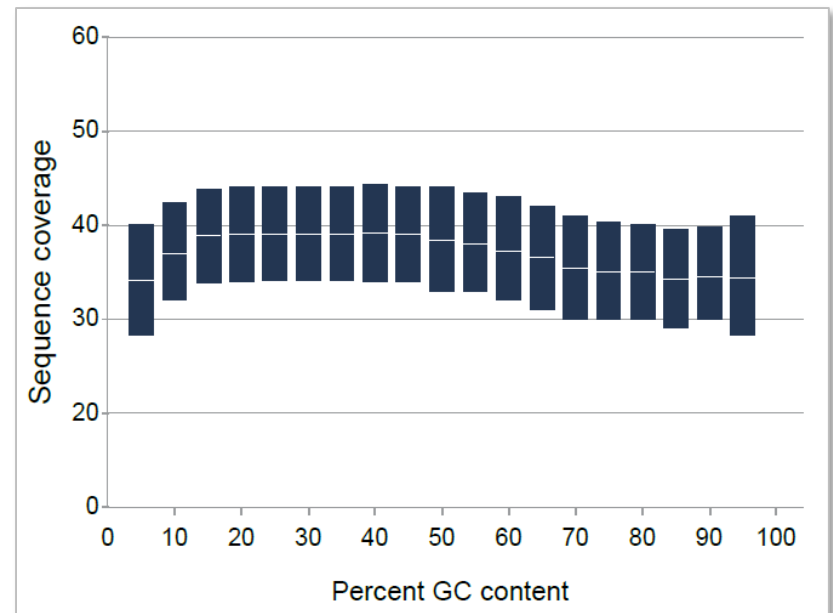
SEQUENCING PERFORMANCE (SEQUEL SYSTEM 3.0 CHEMISTRY)

Long Reads

High Consensus Accuracy

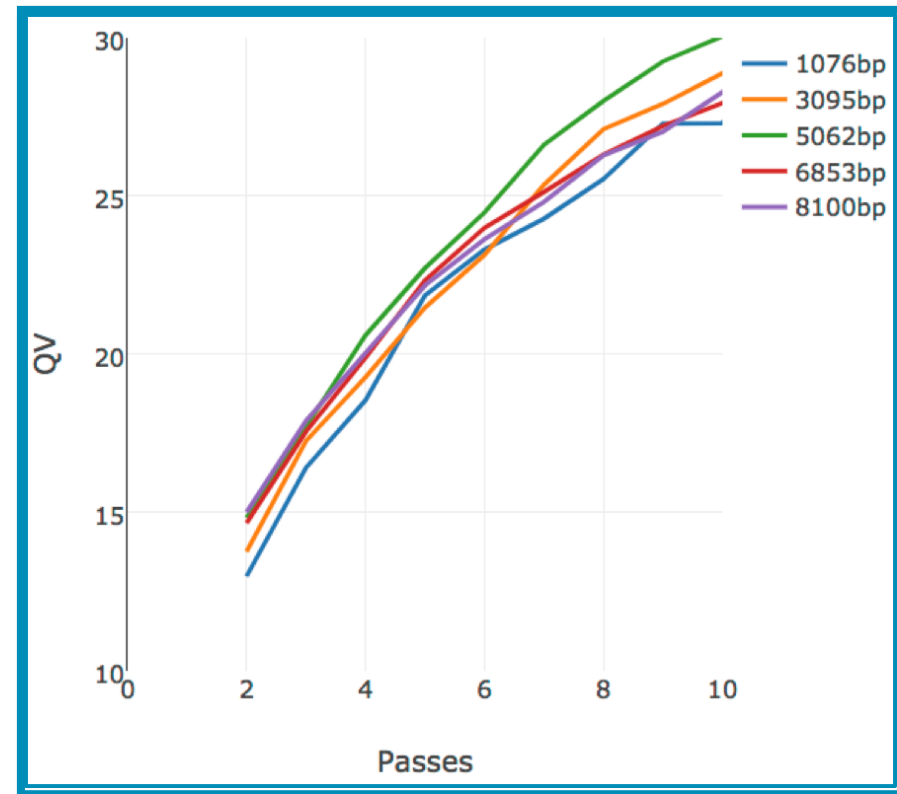
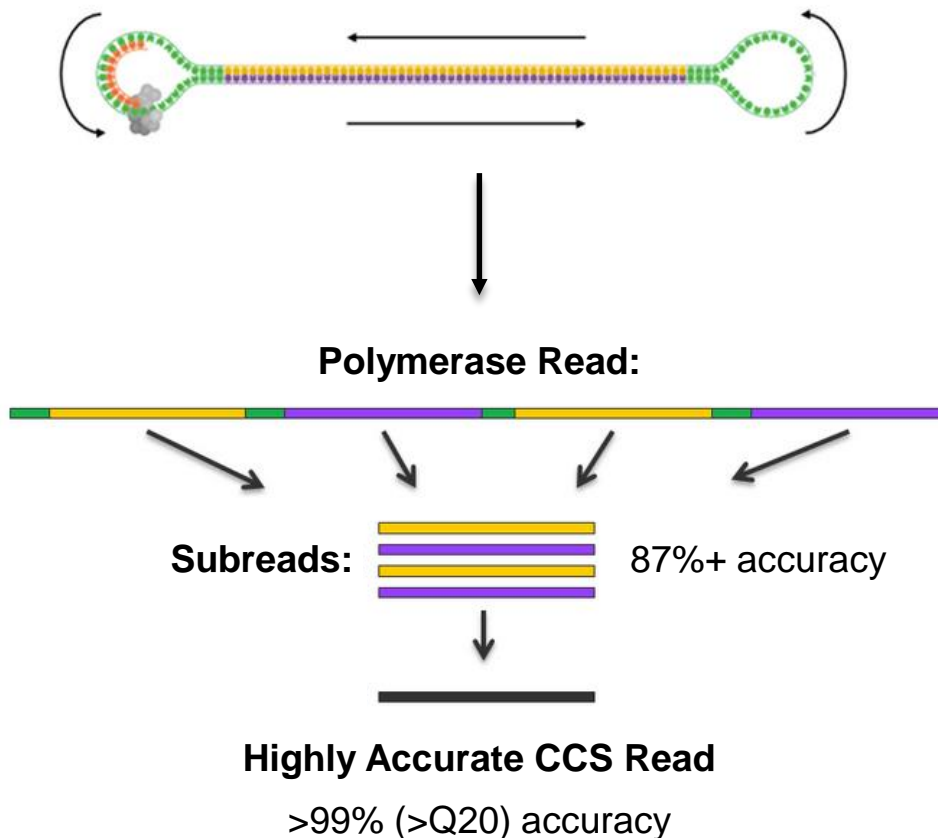
Uniform, Unbiased Coverage

- Minimal GC% or sequence complexity bias



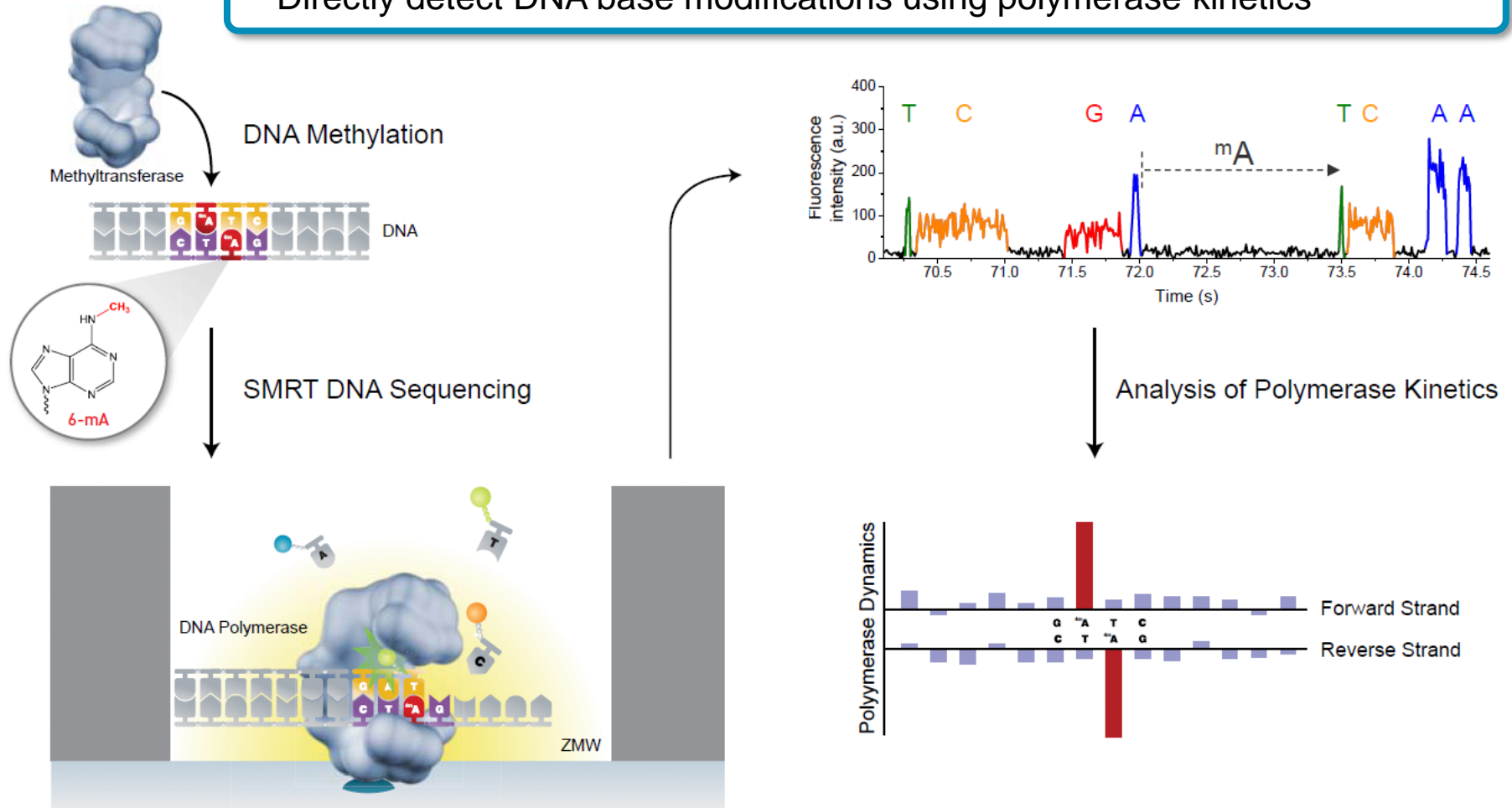
HIGHLY ACCURATE, SINGLE-MOLECULE LONG READS

Significant increase in polymerase read length in the Sequel System 6.0 release increases the number of high fidelity (HiFi), long reads (>Q20 single-molecule read accuracy) for insert sizes up to ~15-20 kb



SIMULTANEOUS EPIGENETIC CHARACTERIZATION

- Directly detect DNA base modifications using polymerase kinetics



DNA polymerization rate is slowed when the polymerase encounters a modified base in the template. Detection of this slowed incorporation rate can be used to infer the presence of bases in the template other than A, C, T or G. This information is automatically generated and processed during every run.



PACIFIC
BIOSCIENCES®



PacBio Long-Read Sequencing Applications

Sequence with Confidence. Pacific Biosciences long-read sequencing provides the most comprehensive view of genomes, transcriptomes, and epigenomes.

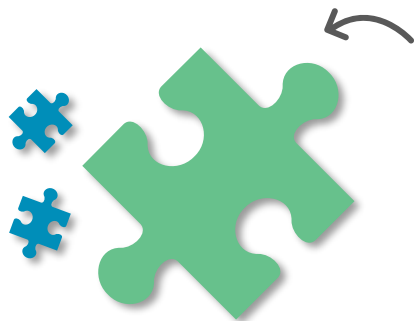


Whole Genome Sequencing for *De Novo* Assembly and Structural Variation Detection

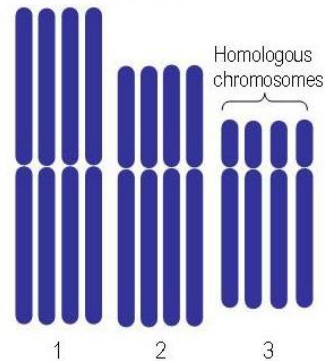
Bring the “W” back to whole genome sequencing

GENERATE GOLD-STANDARD REFERENCE GENOMES AND UNCOVER THE MOST COMPLETE VIEW OF GENOMIC VARIATION WITH PACBIO LONG-READ SEQUENCING

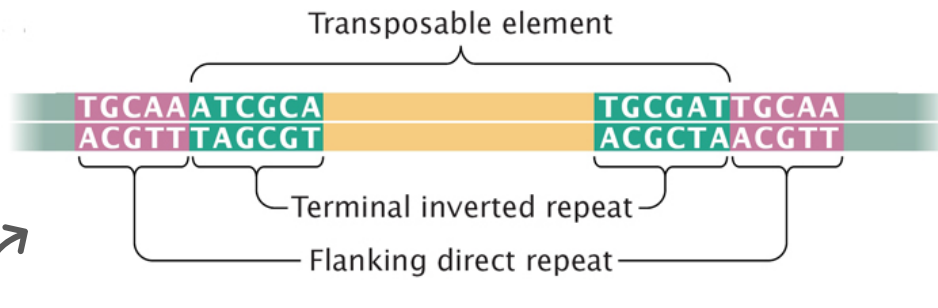
Why Do Long Reads Matter?



bigger is better –
when it comes to
easy assembly



phase haplotypes
in outbred diploids
and polyploids



long reads span repetitive elements
and allow assembly of even the
most complex genomes

even really
big genomes

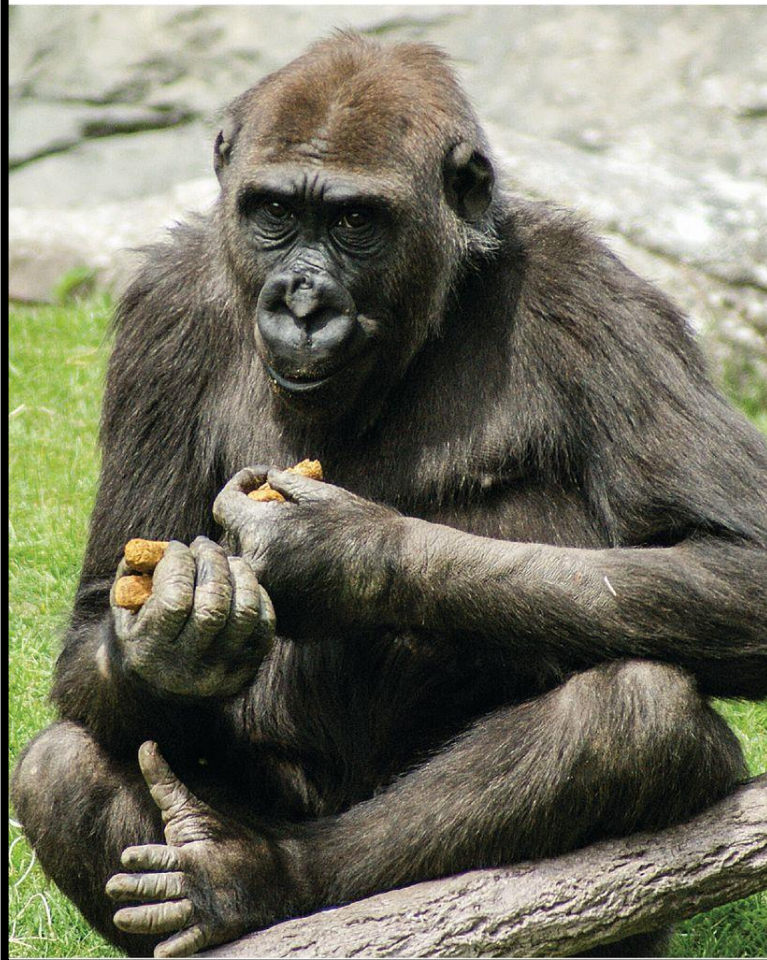


SMRT Sequencing Advantages

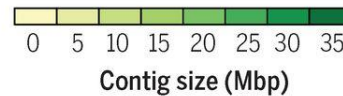
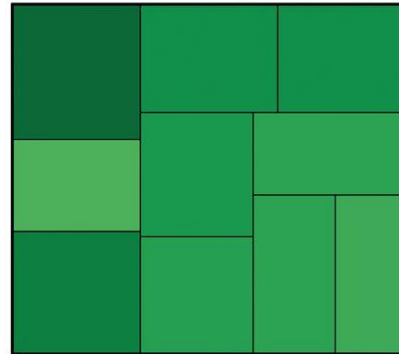


SEQUENCING AND ASSEMBLY OF THE GORILLA GENOME USING LONG-READ VS. SHORT-READ TECHNOLOGY PLATFORMS

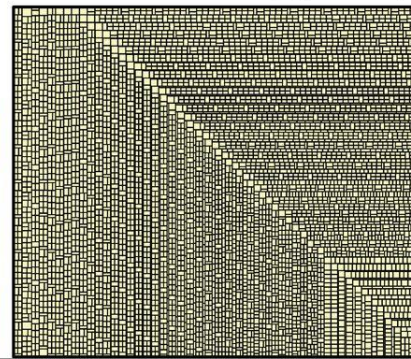
A Susie, reference sample



B Long-read assembly (Susie3)



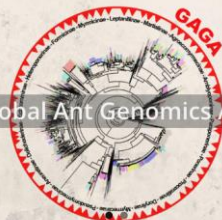
C Short-read assembly (gorGor3)



PacBio long-read assembly showed **>150-fold improvement** over previous short-read assemblies

Treemaps (B and C): Rectangles are the largest contigs that cumulatively make up 300 Mb (~10%) of the assembly, representing the differences in fragmentation of the long-read and short-read gorilla genome assemblies.

The Global Ant Genomics Alliance

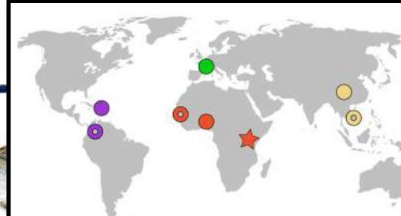


The New York Times

Team of Rival Scientists
Comes Together to Fight Zika



#1MbCtgClub
>350 eukaryotic species!



The MGI Reference Genomes
Improvement Project

WIRED

A COFFEE RENAISSANCE IS
BREWING, AND IT'S ALL
THANKS TO GENETICS



nature
genetics

Putting genome variation
to work in maize
Red leaf pigmentation
Classifying renal cancer



#CanSeq150

nature
biotechnology

EDITORIAL

A reference standard for genome biology



Bloomberg

What the Marijuana Genome Map Means
for the Future of Pot

A breakthrough in genetic research opens the door to more-targeted products
and maybe even pharmaceuticals.

Journal Report
January 17, 2018, 2:00 AM PST



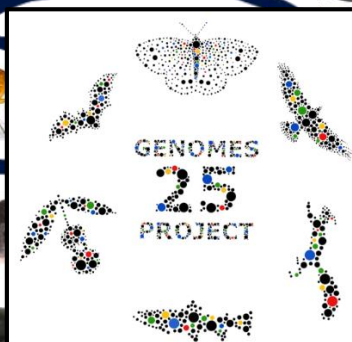
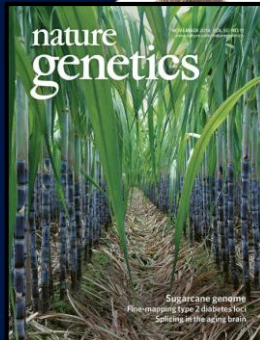
Wine Spectator

Scientists Unravel Cabernet Sauvignon's Genome

Innovative gene-sequencing technology could lead to a better understanding of how
wine grapes evolved and how to adapt them for changing climates

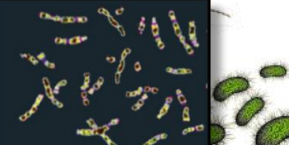


nature
genetics



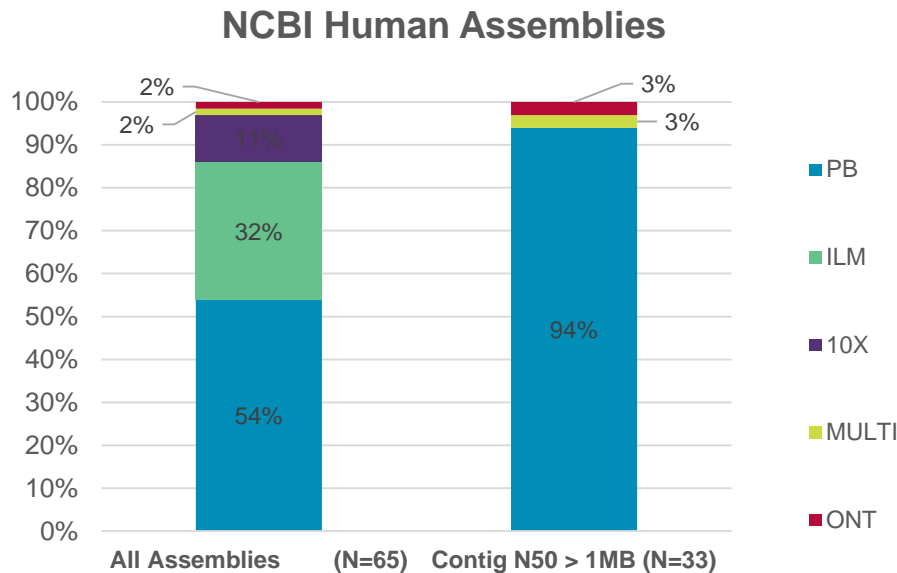
1000 Genomes

A Deep Catalog of Human Genetic Variation



HUMAN WHOLE GENOME SEQUENCING

PacBio is the core technology used for many human reference genome initiatives



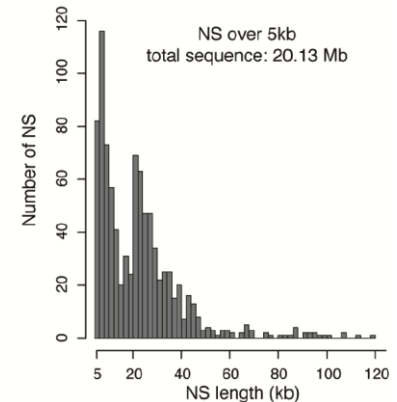
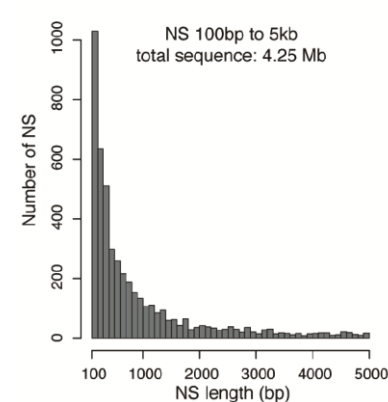
- Generate gold-standard references unique to a population, disease, or individual
- Increase power by matching references to the genetic background of studies
- Access novel types of genetic variation and difficult-to-characterize regions



Article

De Novo Assembly of Two Swedish Genomes Reveals Missing Segments from the Human GRCh38 Reference and Improves Variant Calling of Population-Scale Sequencing Data

Adam Ameur^{1,*}, Huiwen Che¹, Marcel Martin², Ignas Bunikis¹, Johan Dahlberg³, Ida Höijer¹, Susana Häggqvist¹, Francesco Vezzi², Jessica Nordlund³, Pall Olason⁴, Lars Feuk¹ and Ulf Gyllenstein¹



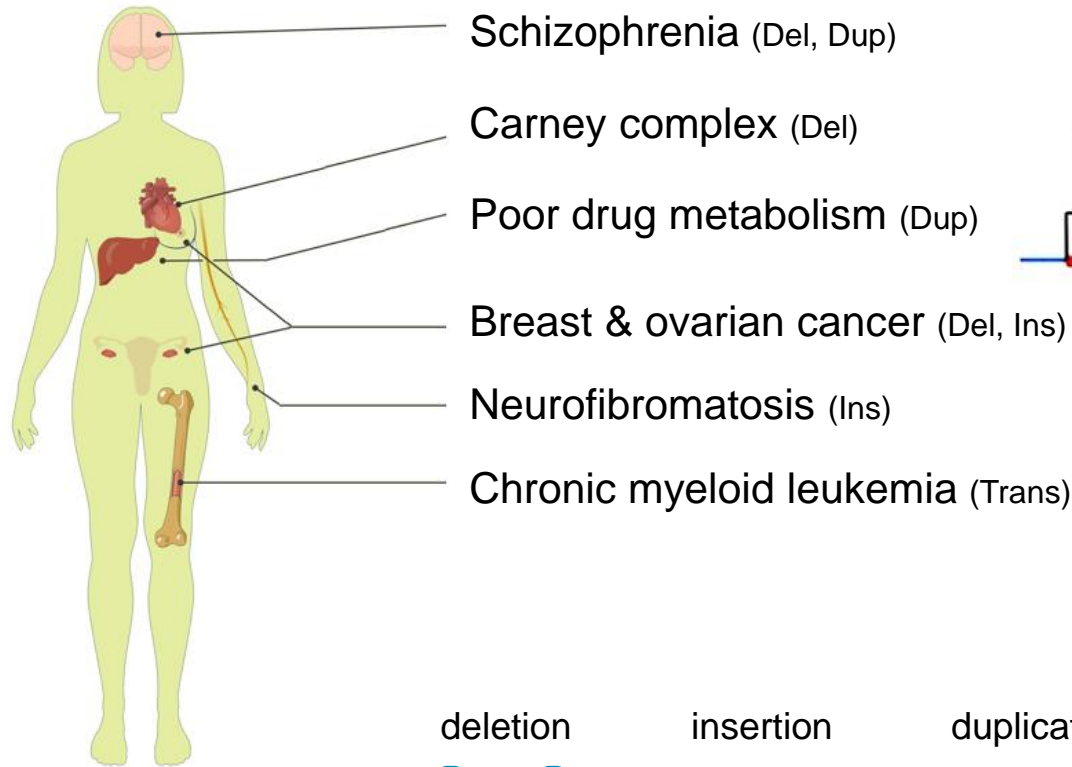
Novel Sequence (NS)

- **~100 Mb** more than Danish short-read assemblies
- **~24 Mb** not found in GRCh38

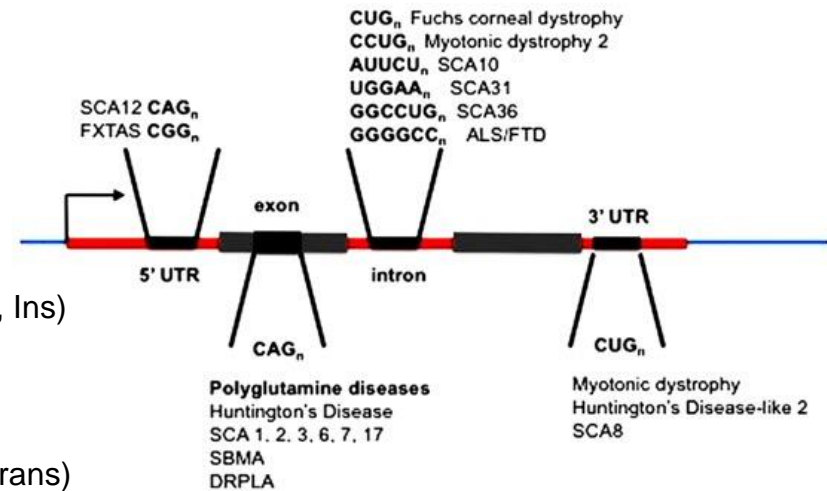
Improved Variant Calls

- **26,000 false positive SNVs eliminated** with improved reference (hg38 + NS)

STRUCTURAL VARIANTS ARE KNOWN TO CAUSE GENETIC DISEASE



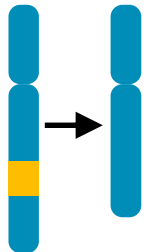
Repeat Expansion Disorders



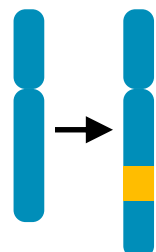
STRUCTURAL VARIANTS

≥ 50 BASE PAIRS / VARIANT

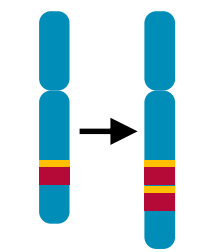
deletion



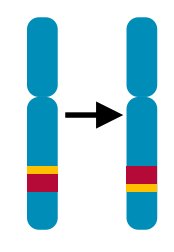
insertion



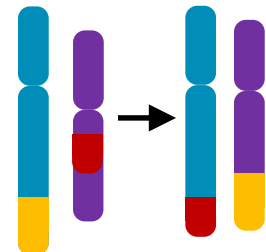
duplication



inversion



translocation

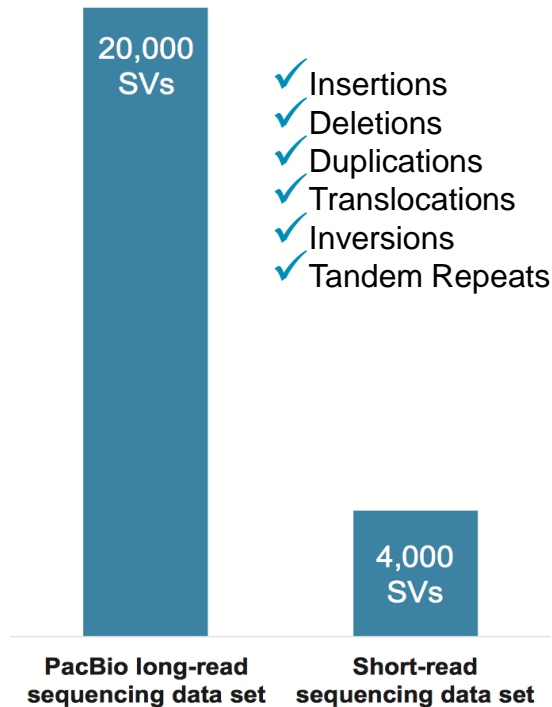


LOW-FOLD PACBIO WGS FOR STRUCTURAL VARIATION DETECTION

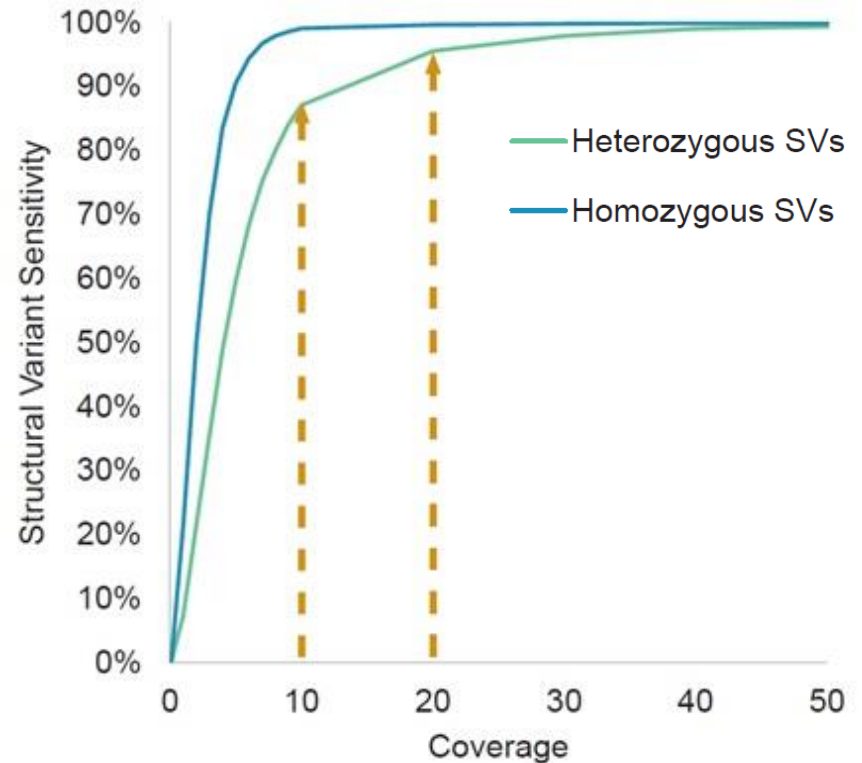
Structural variation accounts for most of the variant bases in the human genome



LONG-READ SMRT SEQUENCING PROVIDES HIGHER SENSITIVITY FOR SV DISCOVERY

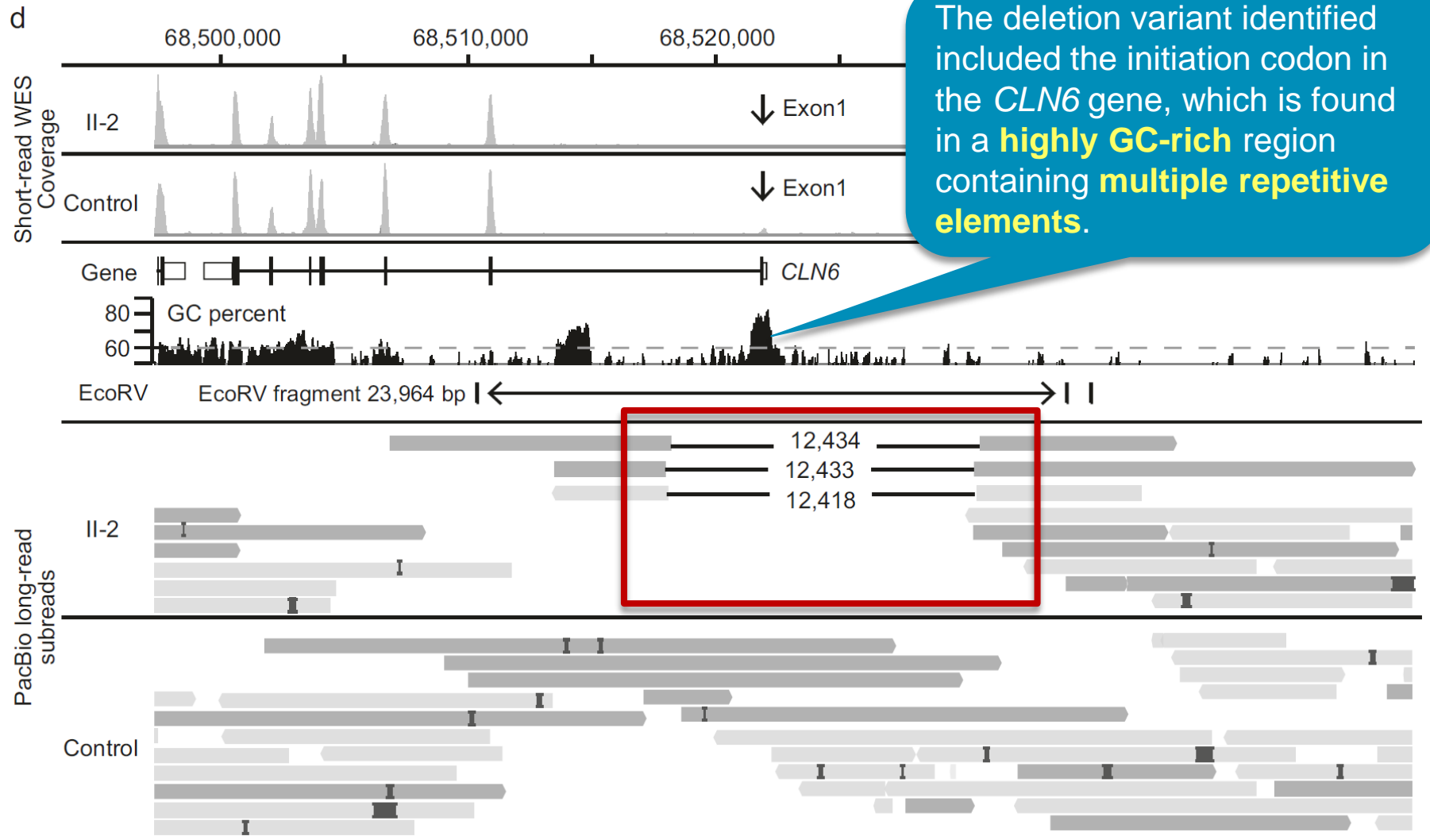


- Low-fold WGS offers a cost-effective option for SV discovery



- Low-fold coverage with PacBio sensitively recalls SVs from a high-coverage HG007233 reference call set.
- Recall is 87% and 95% at 10- and 20-fold coverage, respectively.**

IDENTIFICATION OF A NOVEL 12-KB STRUCTURAL VARIATION IN PROGRESSIVE MYOCLONIC EPILEPSY BY LOW-FOLD (6X) PACBIO WHOLE GENOME SEQUENCING



PLANT AND ANIMAL WHOLE GENOME SEQUENCING

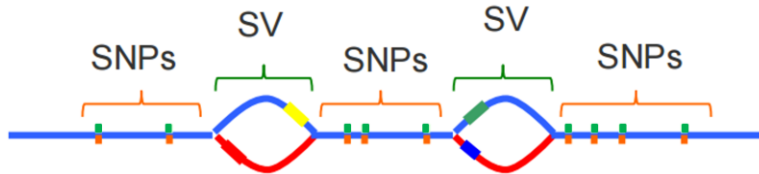
Build better genomes – enable breakthrough discovery

nature **methods**

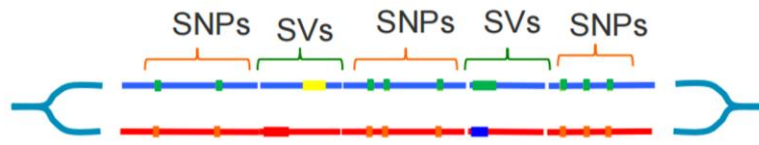
Phased diploid genome assembly with single-molecule real-time sequencing

Chen-Shan Chin^{1,10}, Paul Peluso^{1,10}, Fritz J Sedlazeck², Maria Nattestad³, Gregory T Concepcion¹, Alicia Clum⁴, Christopher Dunn¹, Ronan O'Malley⁵, Rosa Figueroa-Balderas⁶, Abraham Morales-Cruz⁶, Grant R Cramer⁷, Massimo Delledonne⁸, Chongyuan Luo⁵, Joseph R Ecker⁵, Dario Cantu⁶, David R Rank¹ & Michael C Schatz^{2,3,9}

Initial Assembly Graph



Haplotype-Resolved Assembly Graph



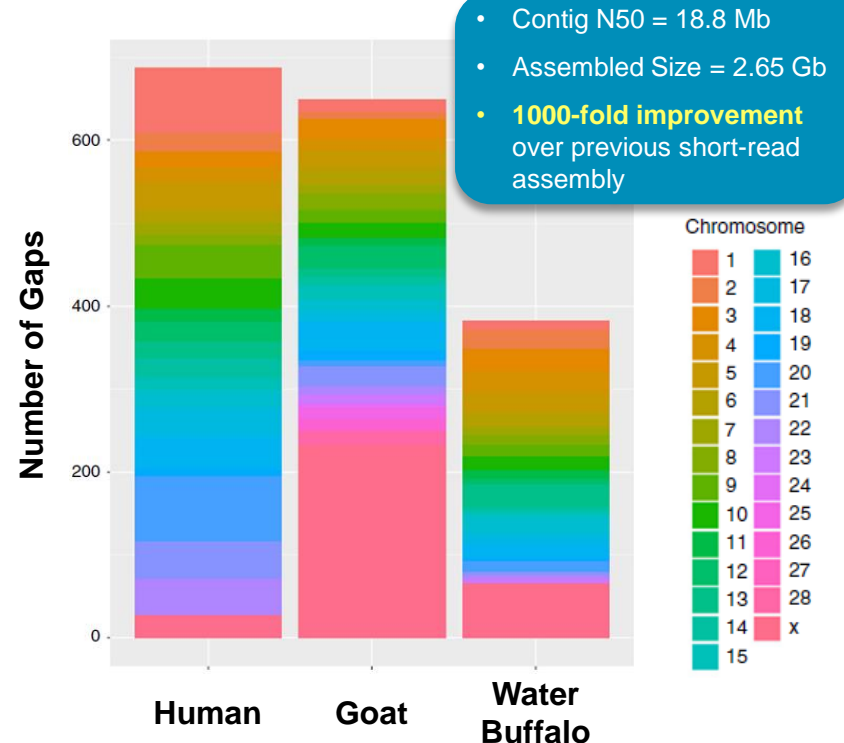
Primary Contig



nature
COMMUNICATIONS

Chromosome-level assembly of the water buffalo genome surpasses human and goat genomes in sequence contiguity

Wai Yee Low¹, Rick Tearle¹, Derek M. Bickhart², Benjamin D. Rosen³, Sarah B. Kingan⁴, Thomas Swale⁵, Françoise Thibaud-Nissen⁶, Terence D. Murphy⁶, Rachel Young⁷, Lucas Lefevre⁷, David A. Hume⁸, Andrew Collins⁹, Paolo Ajmone-Marsan¹⁰, Timothy P.L. Smith¹¹ & John L. Williams¹



LOW DNA INPUT WORKFLOW FOR *DE NOVO* GENOME ASSEMBLY

New low-DNA input protocol puts PacBio-based assemblies in reach for small highly heterozygous organisms that comprise much of the diversity of life






Article

A High-Quality *De novo* Genome Assembly from a Single Mosquito Using PacBio Sequencing

Sarah B. Kingan ^{1,†}, Haynes Heaton ^{2,†}, Juliana Cudini ², Christine C. Lambert ¹, Primo Baybayan ¹, Brendan D. Galvin ¹, Richard Durbin ³, Jonas Korlach ^{1,*} and Mara K. N. Lawnczak ^{2,*}

- ~100 ng of mosquito gDNA used for low-input DNA library preparation protocol (~4-hours)*
- Able to resolve maternal and paternal haplotypes for over 1/3 of the genome
- The method can be applied to samples with starting DNA amounts as low as **≥150 ng per 300 Mb** of genome size

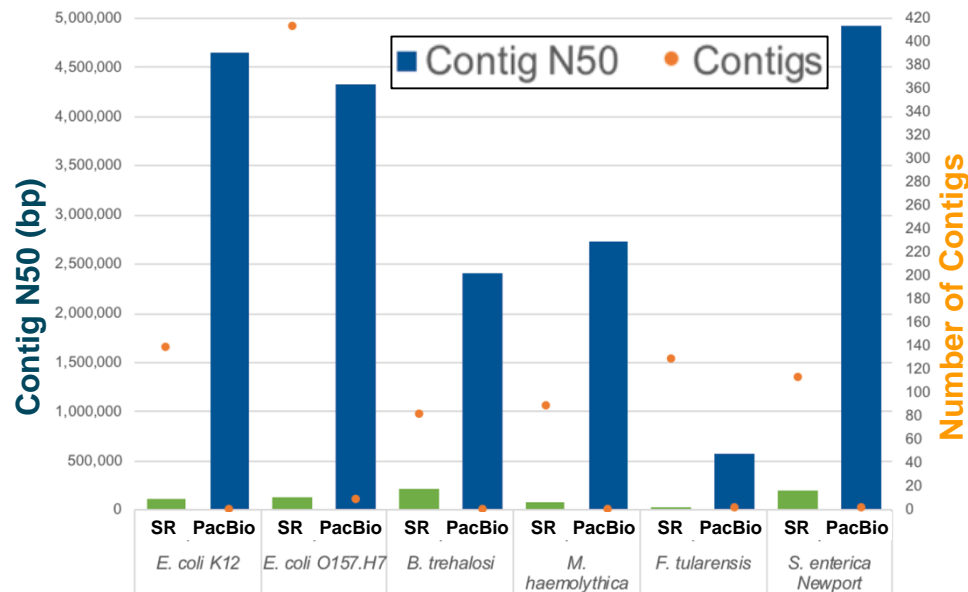
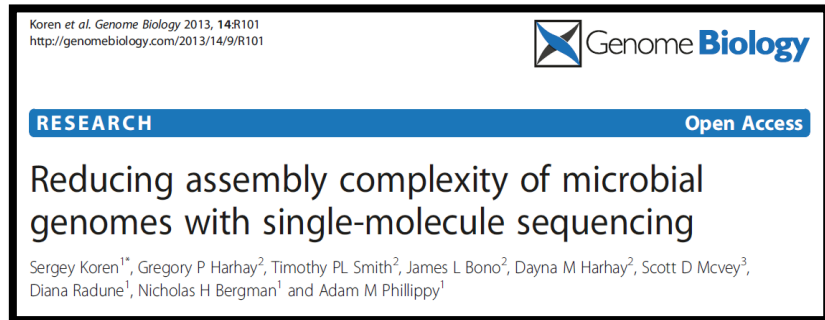
| | | PacBio Assembly | Sanger Assembly |
|-------------------------|-----------------|-----------------|-----------------|
| Primary Contig Assembly | Size (Mb) | 251 | 224 |
| | No. Contigs | 206 | 27,063 |
| | Contig N50 (Mb) | 3.47 | 0.025 |
| Alternate Haplotigs | Size (Mb) | 89.2 | Unresolved |
| | No. Contigs | 830 | N/A |
| | Contig N50 (Mb) | 0.199 | N/A |

“The resulting curated assembly had **high contiguity** (contig N50 = 3.5 Mb) and **completeness** (>98% of conserved genes were present and full-length).”

“In addition, this single-insect assembly now **places 667 (>90%) of formerly unplaced genes into their appropriate chromosomal contexts.**”

MICROBIAL WHOLE GENOME SEQUENCING

Move beyond draft genomes and obtain complete microbial genomes with ease and confidence



Published comparison between different sequencing platforms on the continuity and correctness of genome assemblies for different microbial strains. All PacBio assemblies had QV scores of >60, while the short-read (SR) assemblies had an average QV of 51.3

Koren S. et. al. (2013) *Genome Biology*. 14, R101.

- PacBio microbial genome assemblies are the **gold standard** for both completeness and accuracy
- Only highly accurate, complete genomes **reveal both the SNPs and structural variants** that contribute to drug resistance, virulence, and metabolic evolution
- Only long reads **can resolve repetitive regions encoding important biology**, including synthetic gene clusters, IS elements, active transposons, and phage insertions

- **Affordably assemble** gold-standard genomes by multiplexing up to 16 microbes in one SMRT Cell
- **Identify active RM systems (6mA, 4mC)** directly from whole genome sequencing data
- Study the role of **transposons, phage insertions, and other SVs** in the evolution of virulence
- **Effectively recover plasmids** to track drug resistance and transmission paths

HOSPITAL-ASSOCIATED INFECTIONS

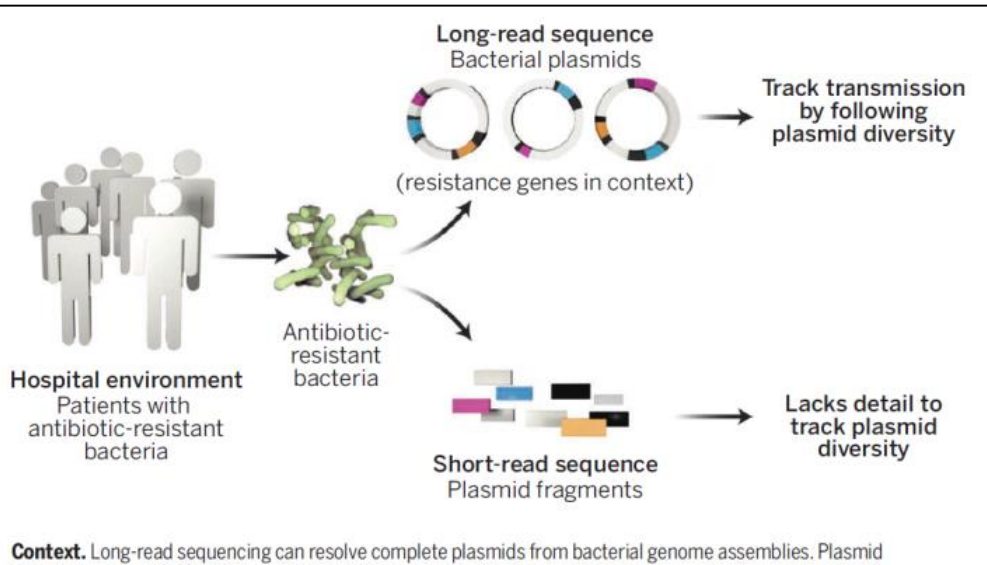
RESEARCH ARTICLE

ANTIBIOTIC RESISTANCE

Single-molecule sequencing to track plasmid diversity of hospital-associated carbapenemase-producing Enterobacteriaceae

Sean Conlan,¹ Pamela J. Thomas,² Clayton Deming,¹ Morgan Park,² Anna F. Lau,³ John P. Dekker,³ Evan S. Snitkin,¹ Tyson A. Clark,⁴ Khai Luong,⁴ Yi Song,⁴ Yu-Chih Tsai,⁴ Matthew Boitano,⁴ Jyoti Dayal,² Shelise Y. Brooks,² Brian Schmidt,² Alice C. Young,² James W. Thomas,² Gerard G. Bouffard,² Robert W. Blakesley,² NISC Comparative Sequencing Program,² James C. Mullikin,² Jonas Korlach,⁴ David K. Henderson,³ Karen M. Frank,^{3*} Tara N. Palmore,^{3*} Julia A. Segre^{1*}

Science
Translational
Medicine
AAAS



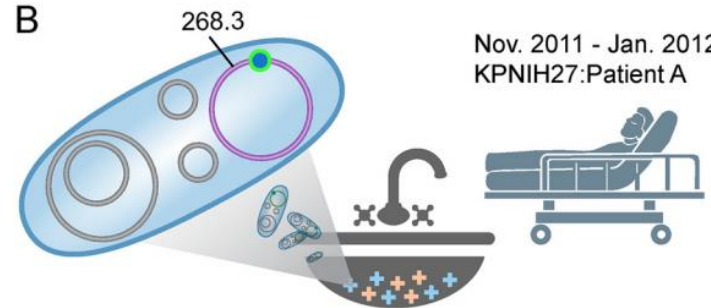
A

Nov. 2010 - Feb. 2011
ECNIH8:Patient Y



B

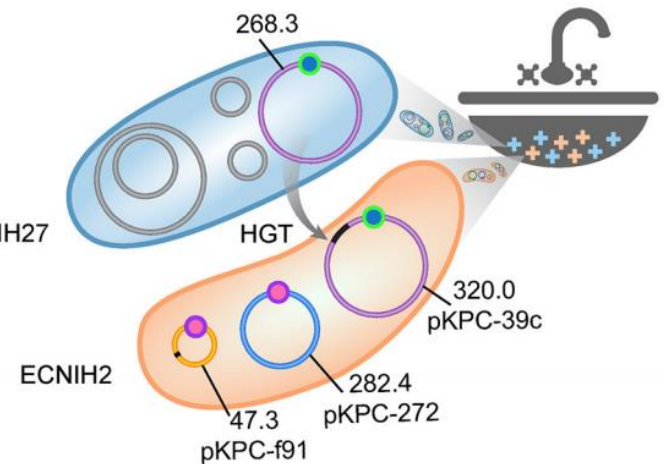
Nov. 2011 - Jan. 2012
KPNIH27:Patient A



C

Jan. 2012
ECNIH2

KPNIH27



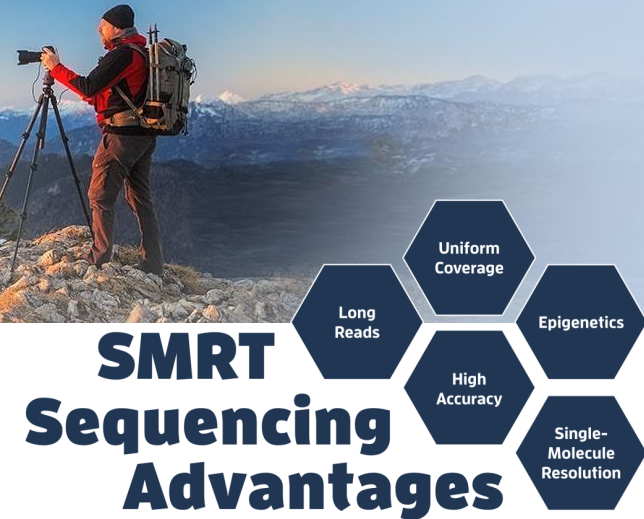
KPC3 Tn4401b

KPC2 IS26-TnpA



Targeted Sequencing

**ACCURATELY DISCOVER AND DETECT ALL
VARIANT TYPES EVEN IN THE HARDEST TO
REACH REGIONS OF THE GENOME**



Benefits of SMRT Sequencing for Amplicons

Simplified Workflow

- Can use a single primer set per target region (up to ~20 kb)

Cost-Effective Multiplexing Options

- Can pool barcoded samples or amplicons from different target regions or projects on a single SMRT Cell and drive down cost/sample

Reduced Time-to-Results

- Less time is required for amplicon design, primer validation, sequencing & assembly, and then redoing preps for dropouts

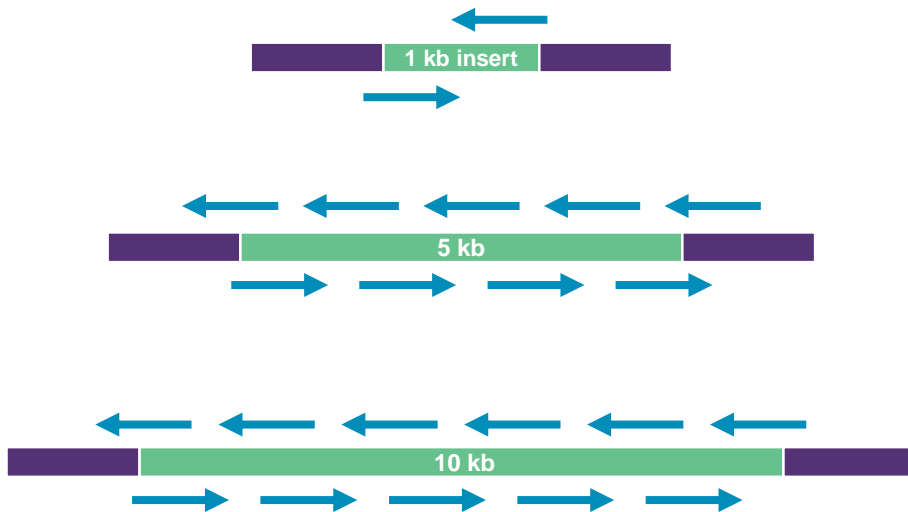
Long Reads with Highest Accuracy and Single-Molecule Resolution

- Enables PacBio to call *all* variants types – single nucleotide variants (SNVs), structural variants (SVs), and copy number variants (CNVs)

TECHNOLOGY COMPARISON – AMPLICON SEQUENCING

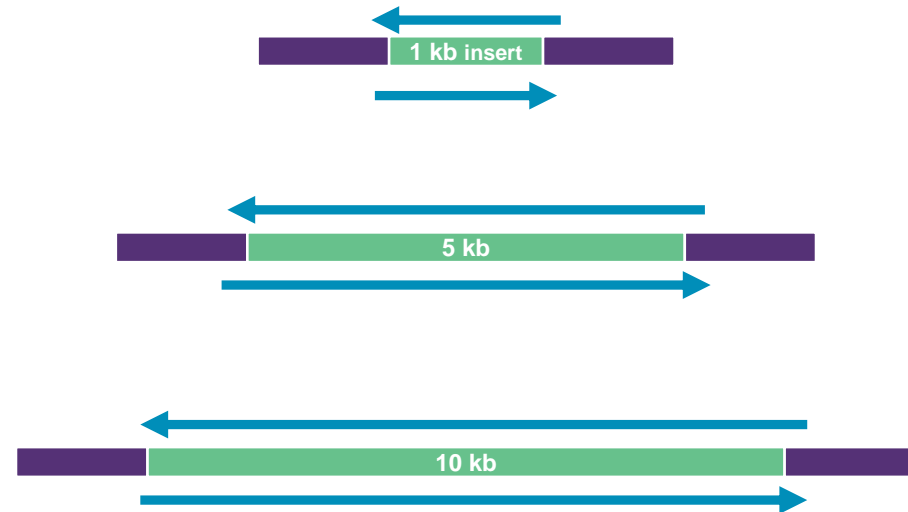
Short-read/Sanger primer design

- Number of primers increases with target length
- Dropouts increase with number of primers



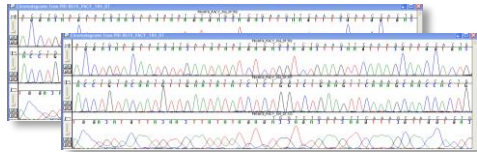
PacBio primer design

- One primer pair for target lengths up to ~10 kb
- Greater flexibility for where to design primers, not necessary to design within target region

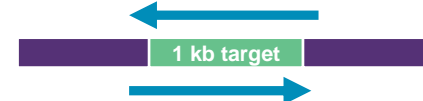


PACBIO AMPLICON SEQUENCING – NO ASSEMBLY REQUIRED

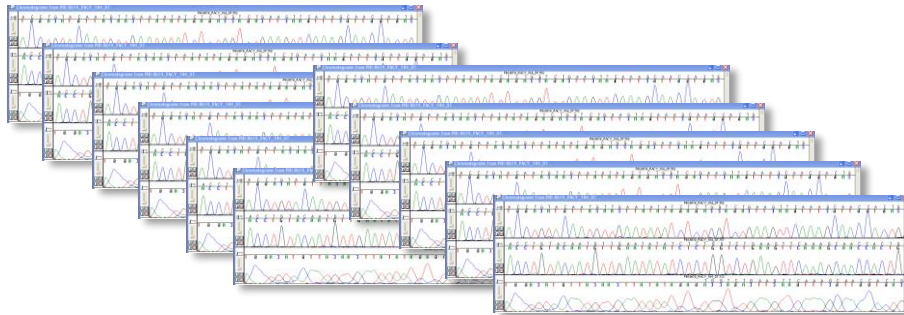
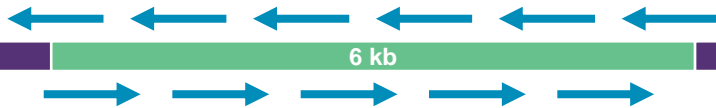
Short-read/Sanger Assembly



PacBio – “No Assembly Required”



```
AGTTGTTAGTCTACGTGGACCGACAAGAACAGTTTCGAATCGGAAGCTTGCTTAACGTAGTCTAACAGT
TTTTTATTAGAGAGCAGATCTCTGATGAACAACCAACGGAAGAGACGGGTCGACCGCTCTTCAATATGC
TGAACACGGCGAGAAACCGCGTGTCAACTGTTTCACAGTTGGCGAAGAGATTCTCAAAAGGATTGCTTTC
AGGCCAAGGACCCATGAAATTGGTGATGGCTTTTATAGCATTCCTTAAGATTCTAGCCATACCTCCAAAC
GCAGGAATTTGGCTAGATGGGGCTCATTCAAGAAAGATGGAGCGATCAAAAGTGTACGGGGTTTCAAGA
AAGAAATCTCAAAACATGTTGAACATAATGAACAGGAGGAAAGATCTGTGACCATGCTCCTCATGCTGCT
GCCACAGCCCTGGCGTTCCATCTGACACCCGAGGGGAGAGCCGCACATGATAGTTAGCAAGCAGGAA
AGAGGAAAATCACTTTTGTAAAGACCTCTGACAGGTGTCAACATGTGCACCCCTTATTGCAATGGATTGG
```



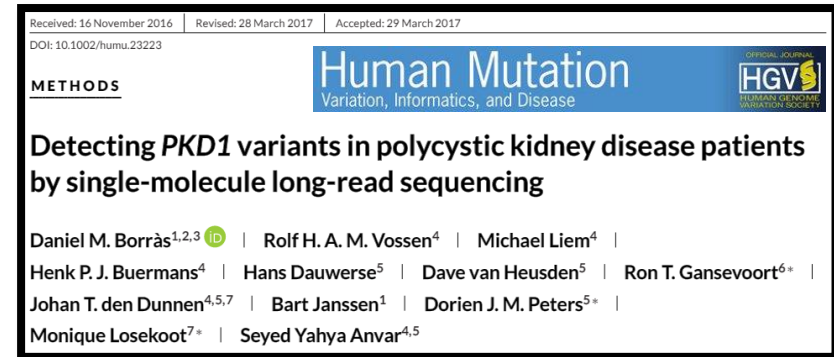
```
AGTTGTTAGTCTACGTGGACCGACAAGAACAGTTTCGAATCGGAAGCTTGCTTAACGTAGTCTAACAGT
TTTTTATTAGAGAGCAGATCTCTGATGAACAACCAACGGAAGAGACGGGTCGACCGCTCTTCAATATGC
TGAACACGGCGAGAAACCGCGTGTCAACTGTTTCACAGTTGGCGAAGAGATTCTCAAAAGGATTGCTTTC
AGGCCAAGGACCCATGAAATTGGTGATGGCTTTTATAGCATTCCTTAAGATTCTAGCCATACCTCCAAAC
GCAGGAATTTGGCTAGATGGGGCTCATTCAAGAAAGATGGAGCGATCAAAAGTGTACGGGGTTTCAAGA
AAGAAATCTCAAAACATGTTGAACATAATGAACAGGAGGAAAGATCTGTGACCATGCTCCTCATGCTGCT
GCCACAGCCCTGGCGTTCCATCTGACACCCGAGGGGAGAGCCGCACATGATAGTTAGCAAGCAGGAA
AGAGGAAAATCACTTTTGTAAAGACCTCTGACAGGTGTCAACATGTGCACCCCTTATTGCAATGGATTGG
GAGAGTTATGTGAGGACACAATGACCTACAAATGCCCGGATCACTGAGACGGAACCAAGATGAGCTTGA
CTGTTGGTGAATGCCACGGAGACATGGGTGACCTATGGAACATGTTCTCAAACTGGTGAACACGGACGA
GACAAACGTTCCGTGCACTGGCAACCAACGTCAGTGGGCTTGGCTAGAAACAAAGCAACGAACTGGATGT
CCTCTGAAGCGGCTTGAACAAATACAAAAGTGGAGACCTGGGCTCTGAGACACCCAGGATTACGGT
GATAGCGCTTTTCTAGCACATGCCATAGGAACATCCATCAACCGAAGAAAGGATCAITTTTATTTCGCTG
ATGCTGGTAATCCATCCATGGCCATGCGGTGCGTGGGAATAGCGCAACAGAGACTTCGTGGAAGGACTGT
CAGGAGCTACGTGGGTGATGTGGTACTGGAGCATGGAAGTTGGCTCACTACCATGGCAAAAGCAAAACC
AACACTGGACATTGAACTCTTGAAGACGGAGGTCAAAAACCTGCGGCTCTGGCGCAAACTGTGCATTGAA
GCTAAAAATATCAAAACACCAACCGATTTCGAGATGTCCACACAAAGGAGAAAGCCAGCTGTGGGAAAGAC
AGGACACGAACTTTTGTGTGTCAGCAGAACGTTCTGGGACAGAGGCTGGGGCAATGGTTGTGGGCTATTCCG
AAAAGGTAGCTTTAATACGTTGTGCTAAGTTTAAAGTGTGTGACAAAACCTGGAAGGAAAGATAGTCCAATAT
GAAAACCTTAAATATTCAGTGAATGTCACCGTACACACTGGAGACCAAGCACCAGTTGGAAATGAGACCA
CAGAACATGGAAACACTGCAACCATTAACACCTCAAGCTCCACGTCGGAATACAGCTGACAGCACTACGG
AGCTCTAACATTGGATTGTTCACTAGAACAGGGCTAGACTTTAATGAGATGTTGTTTGAACAATGAAA
```

MORE EFFICIENT TARGETED SEQUENCING WORKFLOWS ON PACBIO

Example: Detecting *PKD1* variants in polycystic kidney disease patients

Previous Sanger Workflow:

- Four long-range PCR reactions
- >50 nested PCR reactions
- >100 Sanger sequencing reactions
- Multiplex Ligation-dependent Probe Amplification (MLPA)

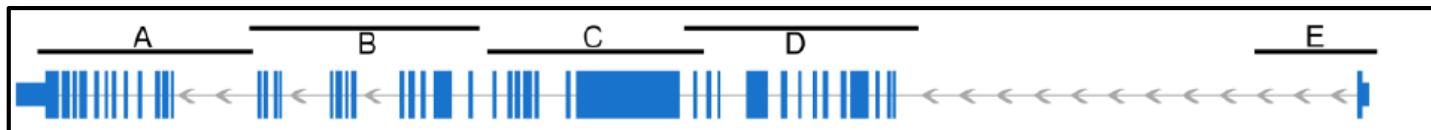


"Our approach provided **high sensitivity** in identifying *PKD1* pathogenic variants, diagnosing 94.7% of the patients."

"We show that **reliable screening of ADPKD patients in a single test without interference of *PKD1* homologous sequences**, commonly introduced by residual amplification of *PKD1* pseudogenes, by direct long-read sequencing is now possible."

New PacBio-based Workflow:

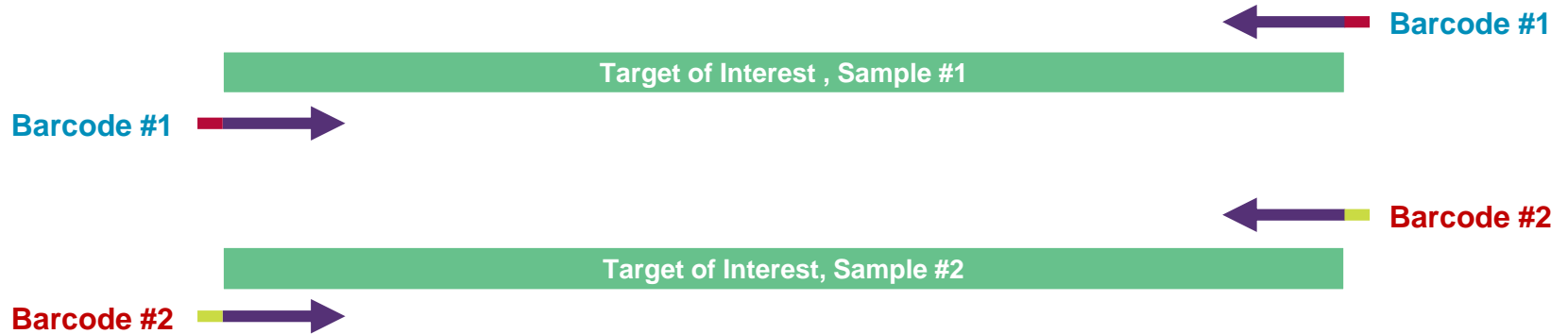
- Five long-range PCR reactions
- Direct full-length amplicon sequencing



PACBIO BARCODES FOR MULTIPLEXED AMPLICON SEQUENCING

A set of 384 barcodes (16 bp length) is optimized for SMRT Sequencing

BARCODED GENE-SPECIFIC PRIMER DESIGN



hundreds of targets
& samples

POOL BARCODED AMPLICONS



1 SMRTBELL LIBRARY PREP



SEQUENCING ON 1 SMRT CELL



ANALYSIS FOR EACH BARCODE



EXAMPLE: POOLING 10,000 BARCODED AMPLICONS ON A SINGLE SMRT CELL ENABLES SIGNIFICANT COST SAVINGS



Asymmetric Barcoding Design:

- 100 forward primer barcodes
 - 100 reverse primer barcodes
 - 800 bp PCR amplicon
- } All possible combinations = 10,000

Sequencing stats:

- Total yield: 6.04 Gb (Chemistry 1.2)
- Total number of reads: 384,274
- Average read length: 15.7 kb

Results:

- All 10,000 bins present with high-quality (>Q40) sequence
- >95% concordance with Sanger data (lower estimate, have found errors in Sanger data)

Hebert et al. BMC Genomics (2018) 19:219
<https://doi.org/10.1186/s12864-018-4611-3>

BMC Genomics

METHODOLOGY ARTICLE

Open Access

A Sequel to Sanger: amplicon sequencing that scales



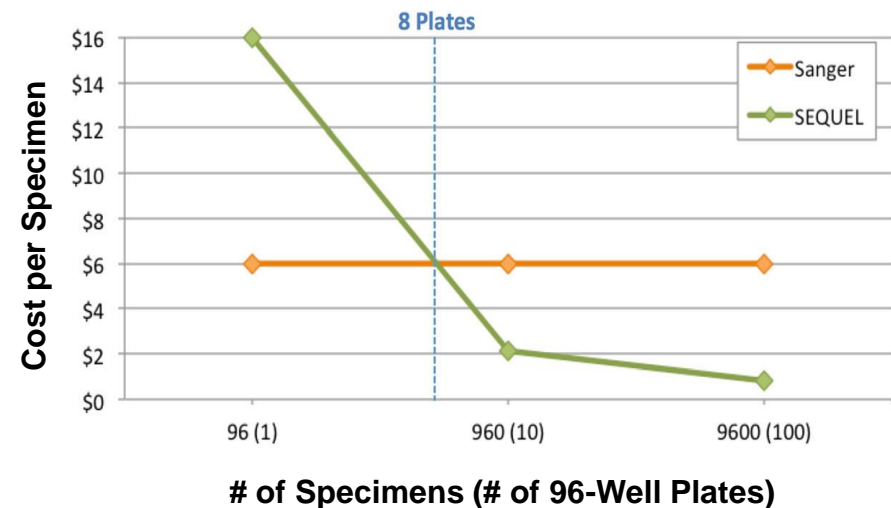
Paul D. N. Hebert^{1*}, Thomas W. A. Braukmann¹, Sean W. J. Prosser¹, Sujeevan Ratnasingham¹, Jeremy R. deWaard¹, Natalia V. Ivanova¹, Daniel H. Janzen², Winnie Hallwachs², Suresh Naik¹, Jayme E. Sones¹ and Evgeny V. Zakharov¹

More accurate:

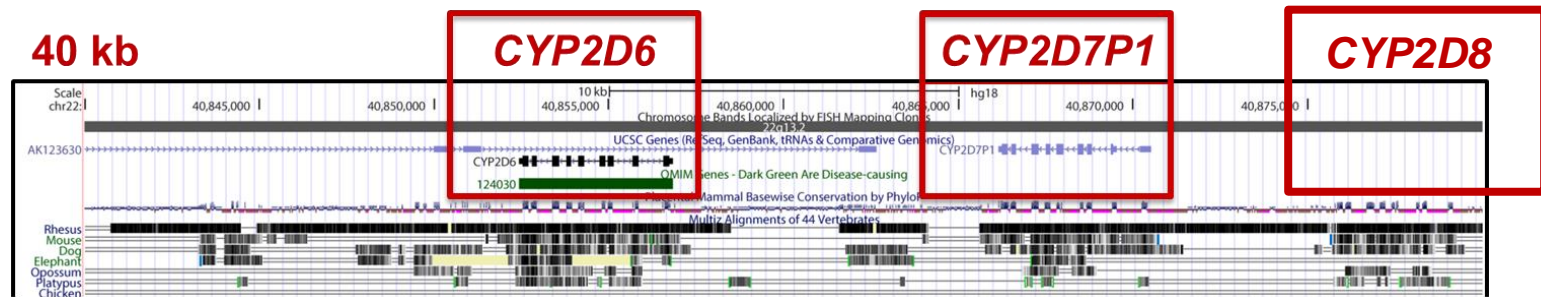
- “SMRT and Sanger sequences were very similar, but SMRT sequencing provided **more complete coverage**, especially for amplicons with homopolymer tracts.”

More cost-effective:

- **40-fold less** compared to Sanger sequencing
- **10-fold less** compared to NGS



FULL-LENGTH SEQUENCING OF 6.6 KB *CYP2D6* AMPLICONS ENABLES PSEUDOGENE DISCRIMINATION



- *CYP2D6* is very difficult to sequence using conventional short-read technology:
 - *CYP2D6* and Pseudogenes *CYP2D7/2D8* are 96% homologous
 - *CYP2D6* locus contains a high frequency of structural deletions and duplications
- Long-read sequencing advantages:
 - Can move away from sequencing only the exonic regions to the entire 6.6-kb *CYP2D6* locus, including the promoter region, all introns, and the downstream regions
 - Enables accurate discrimination of the locus of interest from potential off-target sequences such as pseudogenes

*“For the *CYP2D6* experiments, no high-identity off-target sequence alignments were observed, indicating **pseudogene contamination was not present in our data.**”*

VALIDATION OF A SMRT SEQUENCING PIPELINE FOR HIGH-RESOLUTION ANALYSIS OF *CYP2D6*

- 10 previously characterized DNA samples
- Perfect concordance with previous genotype calls
- One sample showed duplication event; this had been missed by other platforms
- Genotype refined and novel alleles confirmed

| Samples | <i>CYP2D6</i> Diplotype | | TaqMan Copy Number | | <i>CYP2D6</i> SMRT Sequencing | |
|-----------|-------------------------|---------------------|--------------------|--------|-------------------------------|------------------------|
| | Reported ^a | Luminex v3 | Intron 2 | Exon 9 | Upstream / Downstream Copy | Diplotype |
| * NA12244 | *35/*41 | *35/*41 | 2 | 2 | Downstream | *35/*41 ^b |
| NA16688 | *2/*10 | *2/*10 | 3 | 2 | Upstream + Downstream | *2M/*36+*10B |
| NA17222 | *1/*2 | *1/*2 | 2 | 2 | Downstream | *2M/*108 ^c |
| NA17246 | *4/*35 | *4/*35 | 2 | 2 | Downstream | *4/*35 ^{b,d} |
| * NA17247 | *1/*2 | *1/*2 | 2 | 2 | Downstream | *1A/*2M |
| * NA17280 | *2/*3 | *2/*3 | 2 | 2 | Downstream | *3A/*59 ^{d,e} |
| NA17296 | *1/*9 | *1/*9 | 2 | 2 | Downstream | *1A/*9 |
| ASIAN048 | - | *1/*29 ^f | 2 | 2 | Downstream | *1A/*107 ^g |
| HISP291 | - | *1/*17 | 2 | 2 | Downstream | *1A/*17 |
| CAUC053 | - | *1/*6 | 2 | 2 | Downstream | *1A/*6A |

Genotype refinement

Duplication allele characterization

Novel alleles

* validation samples run in triplicate to test intra-run and inter-run reproducibility

CYP2D6 SMRT SEQUENCING AND DIPLTYPE CLARIFICATION

- 14 samples with discrepant results from multiple genotyping platforms
- Provided suballele resolution, genotype refinement, duplicated allele characterization, and discovery of a novel tandem arrangement



Alleles: A, B, C, a, b and c
Genotypes: A/a; B/b and C/c
Haplotypes: ABC and abc
Diplotype: ABC/abc

| Samples | CYP2D6 Diplotype | | TaqMan Copy Number | | CYP2D6 SMRT Sequencing | |
|---------|-----------------------|------------|--------------------|--------|----------------------------|---------------------------|
| | Reported ^a | Luminex v3 | Intron 2 | Exon 9 | Upstream / Downstream Copy | Diplotype |
| NA17289 | *2/*4 | *2/*4 | 2 | 2 | Downstream | *2M/*4 ^b |
| NA17084 | *1/*10 | *1/*10 | 3 | 2 | Upstream + Downstream | *1A/*36+*10B ^c |
| NA17252 | *4/*5 | *4/*5 | 1 | 1 | Downstream | *4/*5 ^b |
| NA17244 | *2A/*4, DUP | *2/*4, DUP | 4 | 4 | Upstream + Downstream | *2Mx2/*4x2 ^b |
| NA17287 | *1/*1(*36/?) | *1/*1 | 2 | 1 | Downstream | *1A/*83 ^d |
| NA09301 | DUP | *1/*2, DUP | 3 | 3 | Upstream + Downstream | *1A/*2x2 ^e |
| NA17218 | *2/*2(*35) | *2/*35 | 2 | 2 | Downstream | *2M/*35 |
| NA17213 | *1/*2(*35) | *1/*35 | 2 | 2 | Downstream | *1A/*35 |
| NA17256 | *2(*35)/*2(*35) | *35/*35 | 2 | 2 | Downstream | *35/*35 |
| NA17243 | *2(*35)/*4 | *4/*35 | 2 | 2 | Upstream + Downstream | *4/*35 ^{b,f} |
| NA17261 | *2(*35)/*4 | *4/*35 | 2 | 2 | Downstream | *4/*35 ^b |
| NA17119 | *1/*2 | *1/*2 | 2 | 2 | Downstream | *1A/*2M |
| CAUC073 | - | ? | 2 | 2 | Downstream | *10B/*109 ^g |
| HISP418 | - | ?, DEL | 2 | 1 | Upstream + Downstream | *5/*36+*41 ^h |

Suballele resolution

Genotype refinement

Duplication allele characterization

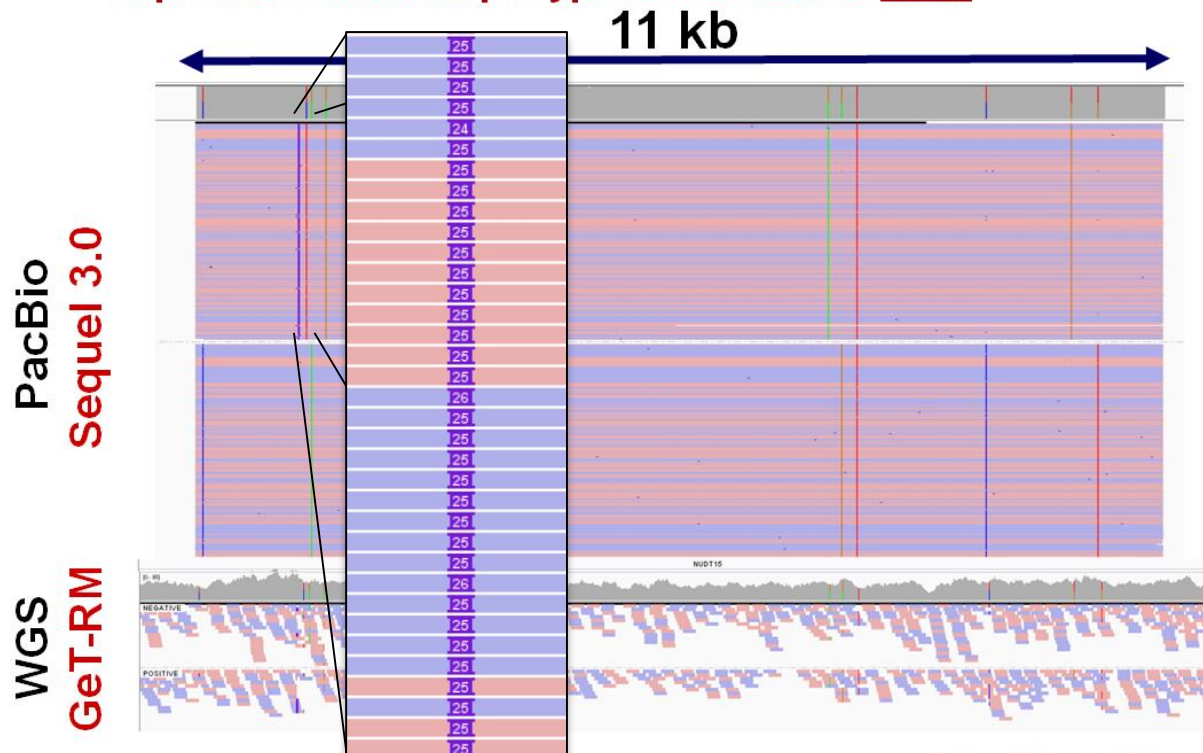
Novel tandem arrangement

11 KB LONG AMPLICON EXAMPLE: PHASING AND DETECTION OF A 25 BP INSERTION MISSED BY SHORT-READ SEQUENCING

PGX: PGx gene SMRT SEQUENCING – Sequel v3.0

2. New v3.0 chemistry: ~11 kb PGx gene amplicon

- **Triplicate variant / haplotype concordance: 100%**

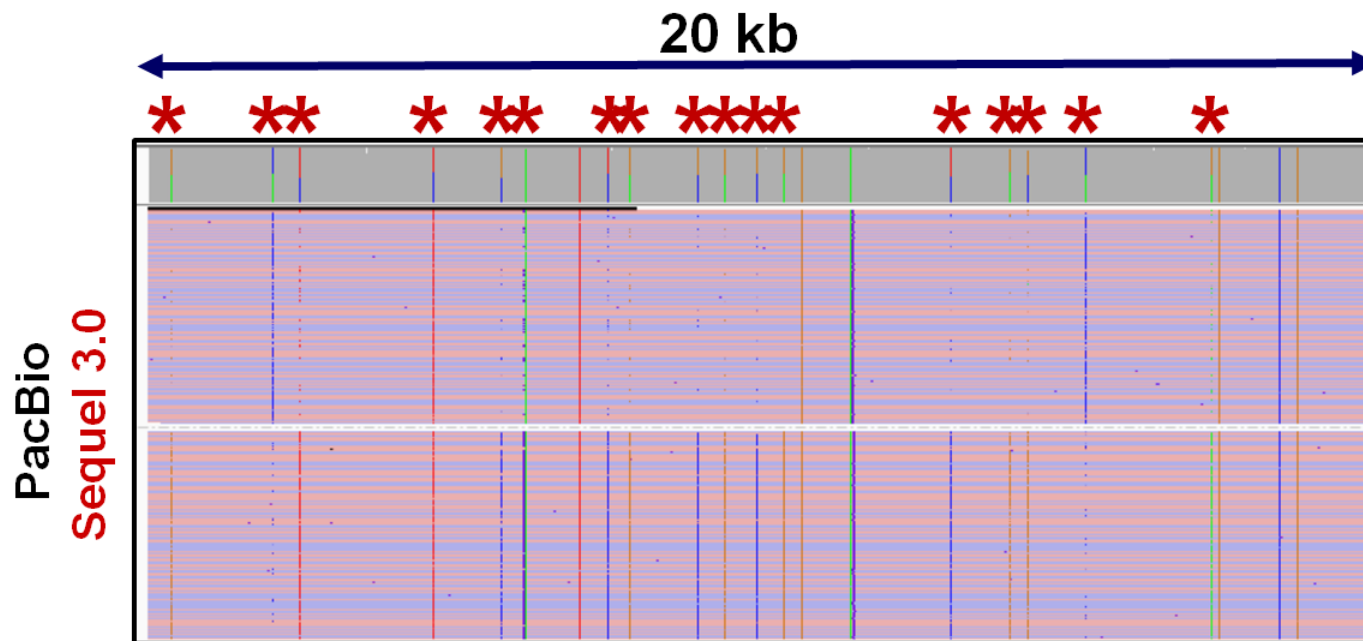


Botton MR, et al. *In preparation.*

LONG-RANGE PCR AMPLICON EXAMPLE: PHASING ACROSS A ~20 KB MENDELIAN DISEASE GENE

CLINICAL: SMRT SEQUENCING – **Sequel v3.0**

3. New v3.0 chemistry: ~20 kb Mendelian disease gene amplicon:



Cody N, et al. *In preparation.*



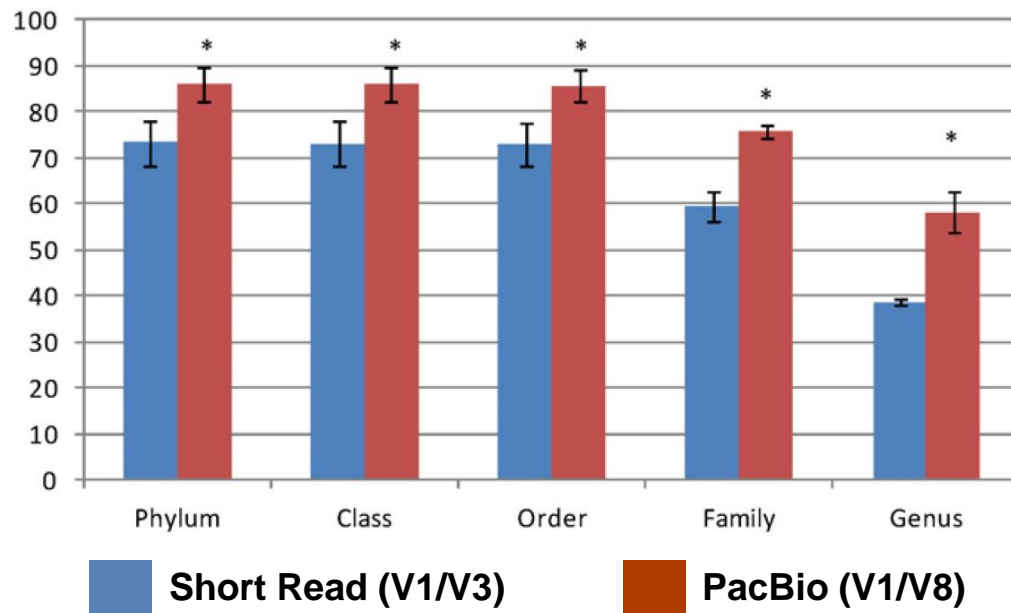
Analysis of Complex Populations

RESOLVE MICROBIAL POPULATION COMPLEXITY WITH ACCURACY AND CONFIDENCE

Obtain species- and strain-level resolution of microbial community form and function with full-length 16S profiling and long-read metagenomic profiling



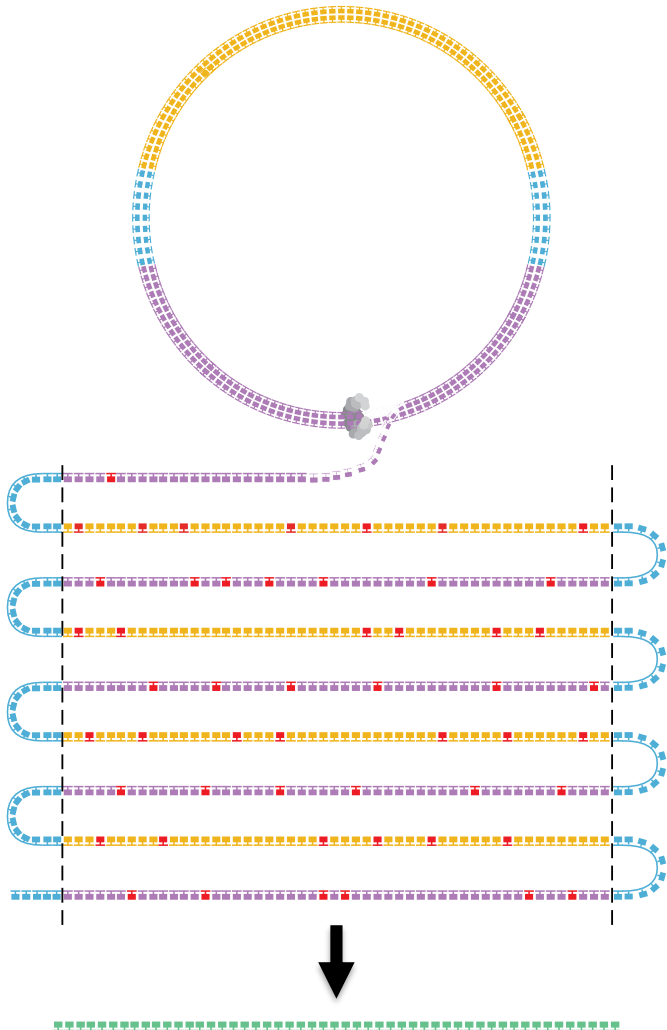
SMRT Sequencing for Metagenomics Pairs Long Reads with High Accuracy



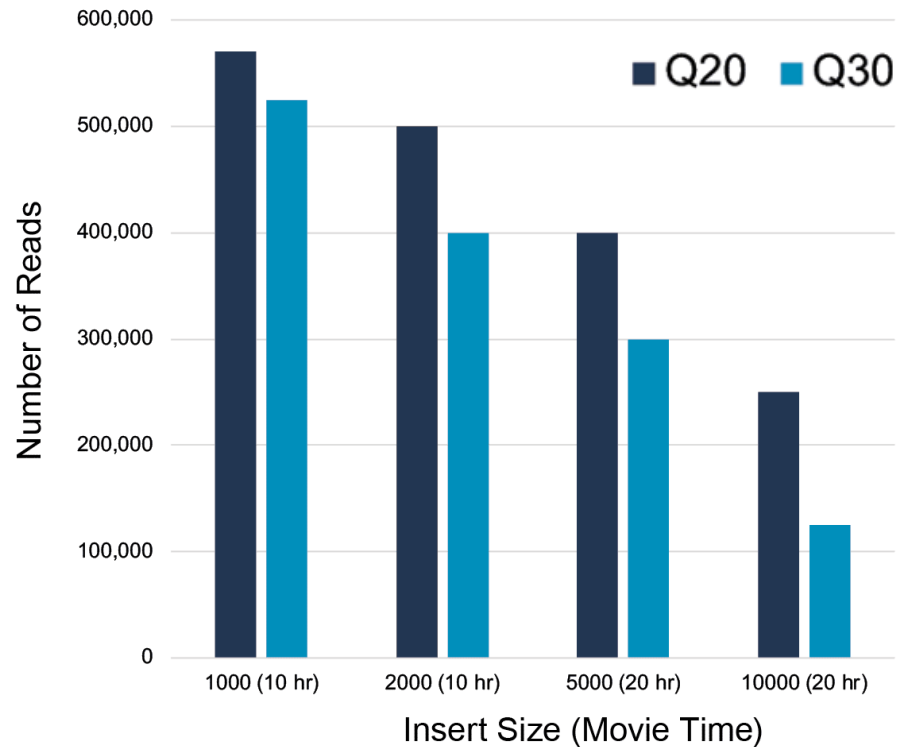
Proportion of 16S rRNA PacBio and short-read generated sequences successfully classified and assigned at five taxonomic levels. Error bars represent standard errors of the mean. *P < 0.05. (P.R. Myer et al. / Journal of Microbiological Methods 127 (2016) 132–140)

- Identify microbial community members with strain-level resolution using **full-length 16S rDNA** sequencing
- Understand key community functions by sequencing complete operons at 99.9% accuracy with **long-insert metagenomic profiling** – no assembly required
- Discover novel genes and gene clusters by reconstructing multi-kilobase long contigs with **whole genome shotgun metagenomic assembly**

PACBIO LONG READ LENGTHS CAN BE USED TO GENERATE HIGH-FIDELITY CCS SEQUENCES OF SINGLE MOLECULES



Highly Accurate CCS Read



- Up to **500,000 Q30 full-length 16S sequences** can be obtained from a single SMRT Cell 1M in a 10-hour run
- For shotgun metagenomics, generate up to **250,000 Q20 reads** from a 10 kb insert library in a 20-hour run

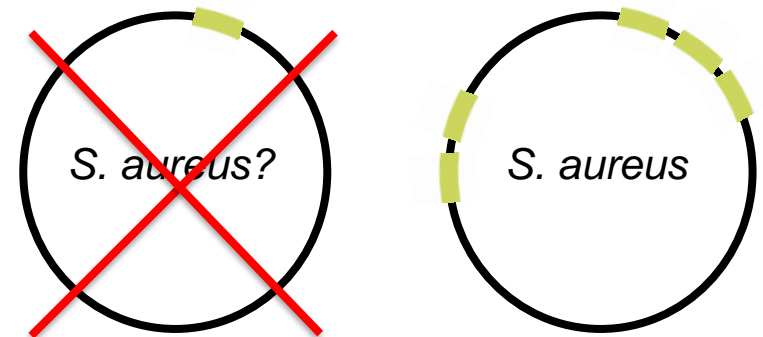
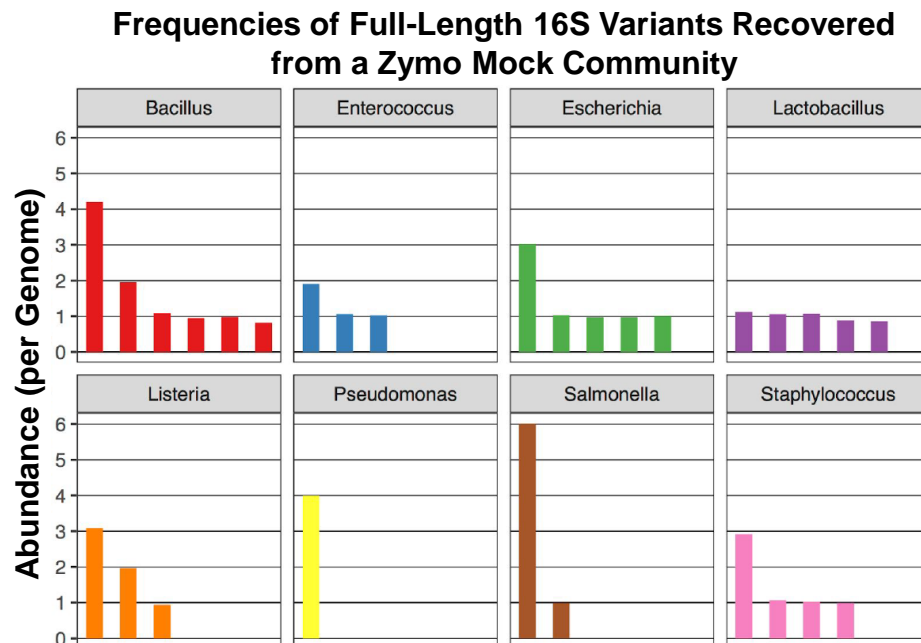
SMRT SEQUENCING REVEALS MULTIPLE DISTINCT 16S SEQUENCES PER BACTERIAL GENOME

- Every bacteria has multiple copies of the 16S housekeeping gene, but in many cases they are not perfect duplicates
- PacBio full-length 16S sequencing revealed multiple distinct 16S variants per bacterial genome
- 16S variants appear in integer ratios that reflected their copy number in each genome

doi: <https://doi.org/10.1101/392332>



Re-analysis of a *Staphylococcus aureus* isolate revealed that previously reported “systematic errors” were actually ground truth



“The differences between these intragenomic variants may have been misinterpreted as systematic errors, perhaps because **the short-read genome assembly that was used as the ground truth contained only one of the five rRNA operons** in the *S. aureus* genome”

ACCURATELY MAPPING MYCOBIOTA ITS1 SEQUENCES TO UNDERSTAND HUMAN HEALTH



Fungal ITS1 Deep-Sequencing Strategies to Reconstruct the Composition of a 26-Species Community and Evaluation of the Gut Mycobiota of Healthy Japanese Individuals

Daisuke Motooka^{1†}, Kosuke Fujimoto^{2,3†}, Reiko Tanaka⁴, Takashi Yaguchi⁴, Kazuyoshi Gotoh^{1,5}, Yuichi Maeda^{2,3}, Yoki Furuta², Takashi Kurakawa², Naohisa Goto¹, Teruo Yasunaga¹, Masashi Narazaki³, Atsushi Kumanogoh³, Toshihiro Horii¹, Tetsuya Iida¹, Kiyoshi Takeda² and Shota Nakamura^{1*}

PacBio CCS Sequencing of PCR amplicons targeting the fungal internal transcribed spacer (ITS) region, ITS1, **most accurately represented the metagenomic population profile** compared to other technologies

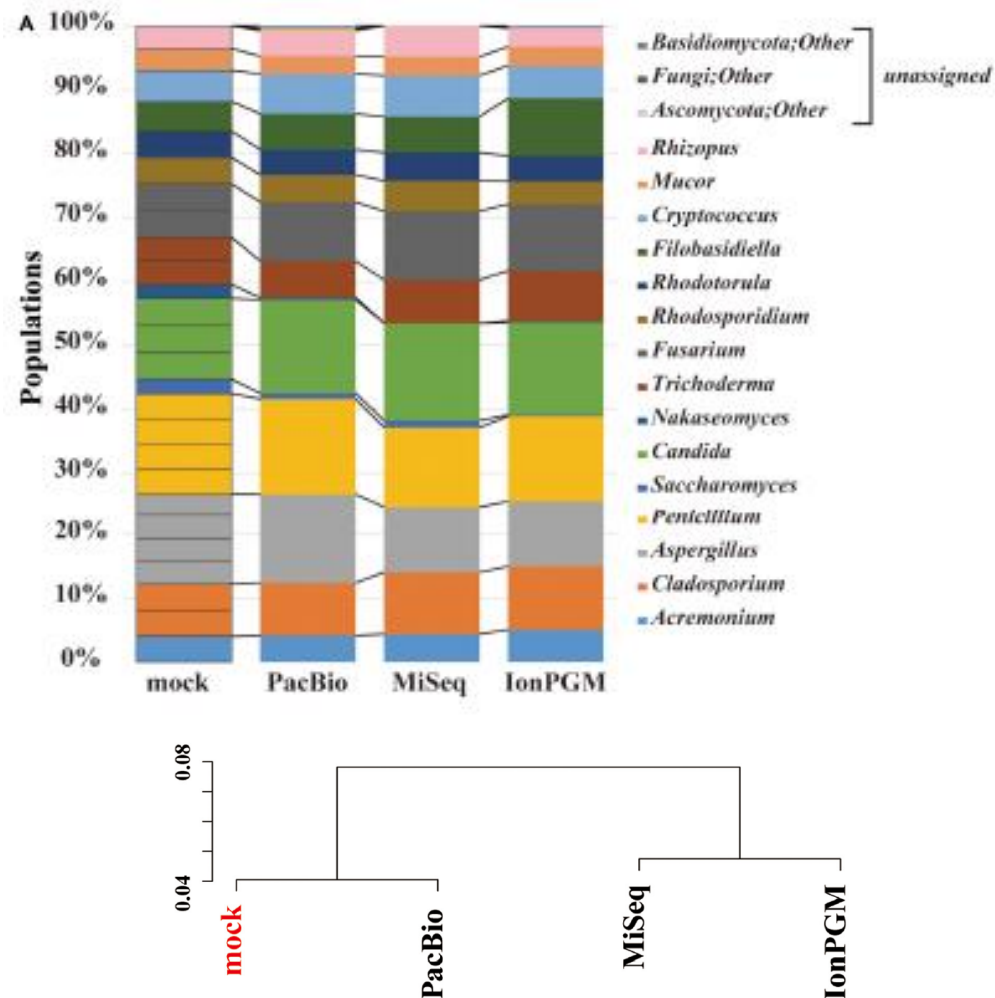
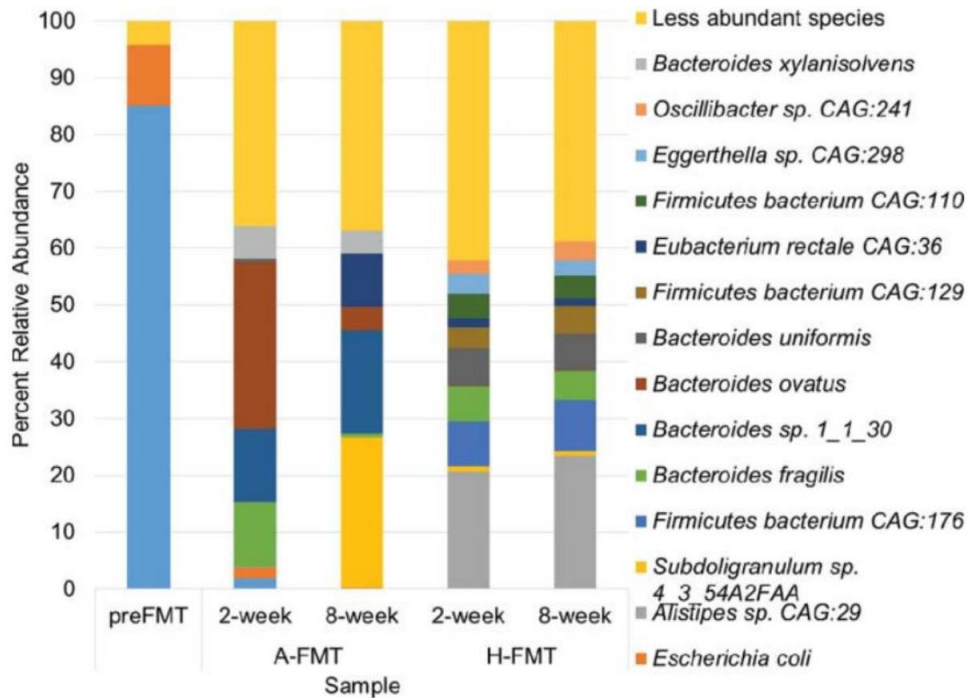


Figure shows the relative abundance of major fungal genera in the mock community.

LONG-READ SHOTGUN METAGENOMIC PROFILING OF MICROBIOME DIVERSITY OF FECAL MICROBIOTA TRANSPLANTATION (FMT) PATIENTS

Long, single-molecule CCS reads generated high-resolution metagenomic profiles to the species and strain level for pre- and post-FMT samples from patients suffering from chronic *C. difficile* infection



Species composition of samples characterized using the PacBio Sequel platform. A-FMT: autologous FMT; H-FMT: heterologous FMT. Resolution of communities at the species level revealed a greater shift in community composition among A-FMT samples than H-FMT samples. A-FMT samples are characterized by fluctuations in abundance of species predominantly within the genus *Bacteroides*. In contrast, H-FMT communities appeared more taxonomically stable and are comprised of a highly abundant species of *Alistipes* and more consistent distribution of *Bacteroides* spp.

Hall, R. et al. (2017). Poster presented at AGBT 2017, Hollywood Beach, Florida.

Comparison of information yield per unit cost for different sequencing and gene-calling methods of identifying full-length proteins from selected spike-in control genomes.

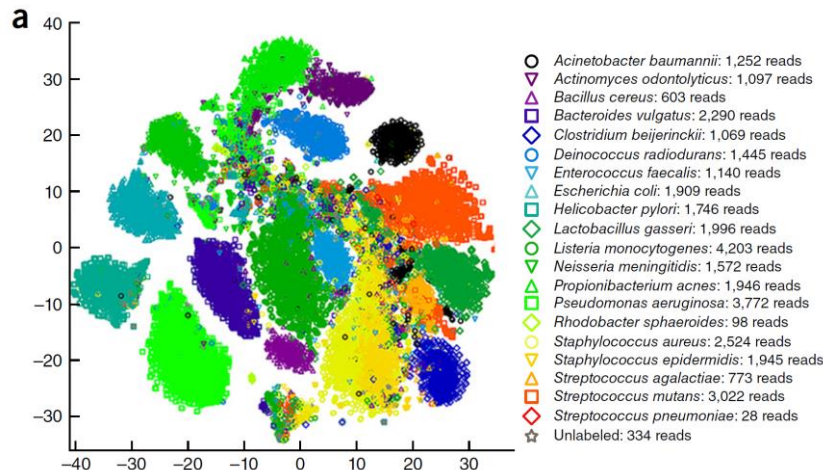
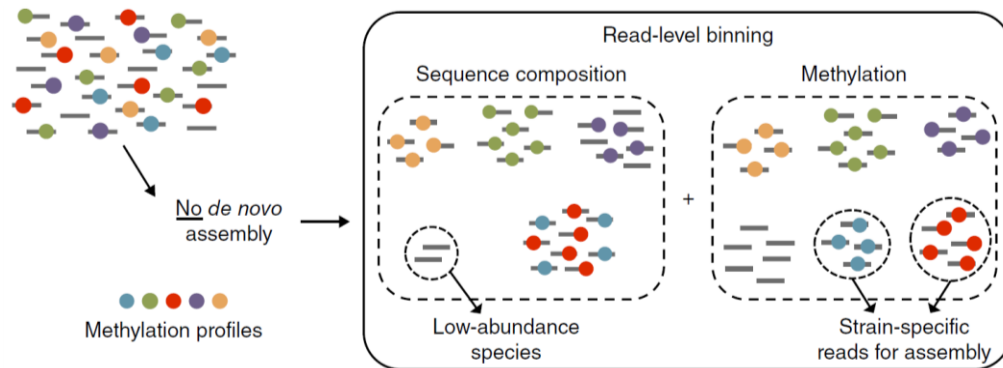
| Organism | Short Read | | PacBio | |
|----------|------------|---------------|----------|---------------|
| | Prodigal | FragGene Scan | Prodigal | FragGene Scan |
| AG | 1,321 | 1,182 | 2,180 | 2,038 |
| BL | 1,191 | 986 | 2,028 | 1,744 |
| PL | 1,197 | 990 | 2,097 | 1,814 |

“By normalizing data to calculate # of unique proteins predicted per \$1,000, researchers at **Second Genome** found that **PacBio sequencing was actually twice as cost-effective as short-read technology at discovering complete genes from the same metagenomic DNA sample.**”

“Whereas short-read technology predicted ~17,000 full-length proteins per \$1,000 of data, PacBio data yielded ~36,000 predicted proteins.”

PacBio BLOG Post Mar 27, 2019. <https://www.pacb.com/blog/>

METAGENOMIC BINNING AND ASSOCIATION OF PLASMIDS WITH BACTERIAL HOST GENOMES USING DNA METHYLATION



Top figure: Overview of metagenomic binning using DNA methylation detected in SMRT long reads. Read-level binning by methylation profiles can segregate reads from multiple strains for the purpose of separate, strain-specific *de novo* genome assemblies.

Bottom figure: t-SNE scatterplot shows binning of assembled contigs and raw reads from a Human Microbiome Project mock community sample using sequence composition and DNA methylation profiles.

**nature
biotechnology**

Article | Published: 11 December 2017

Metagenomic binning and association of plasmids with bacterial host genomes using DNA methylation

John Beaulaurier, Shijia Zhu, Gintaras Deikus, Ilaria Mogno, Xue-Song Zhang, Austin Davis-Richardson, Ronald Canepa, Eric W Triplett, Jeremiah J Faith, Robert Sebra, Eric E Schadt & Gang Fang

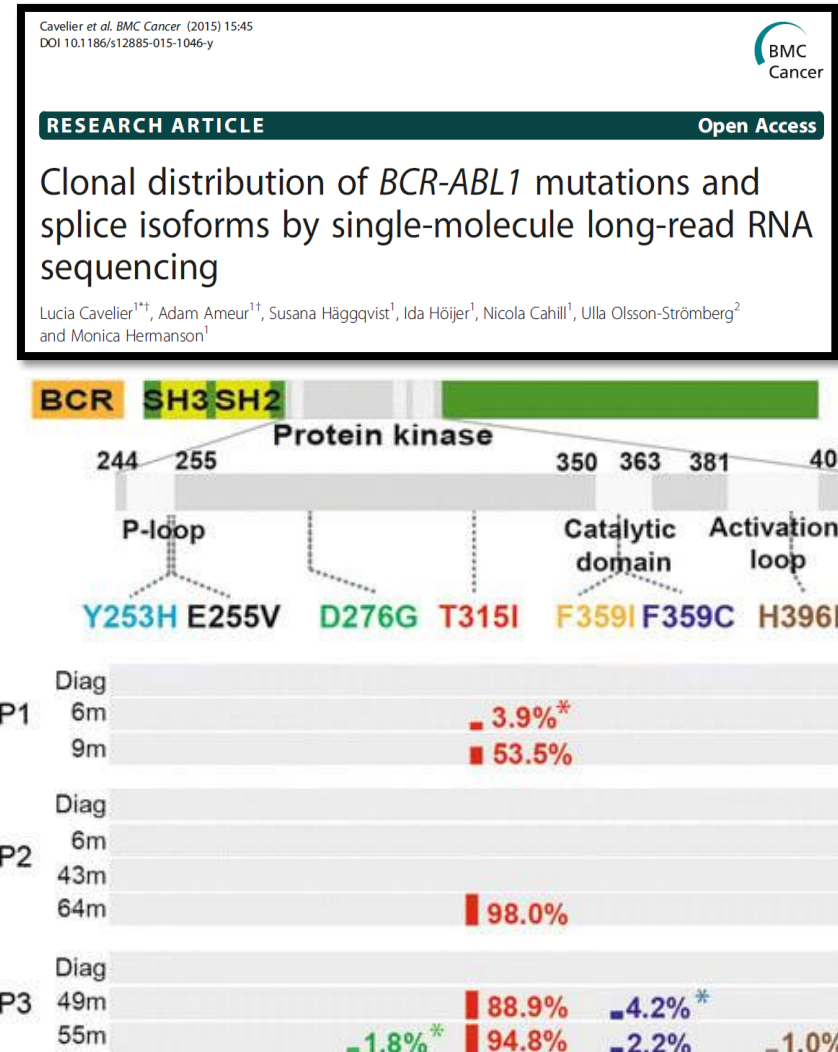
- DNA methylation patterns can be exploited as highly informative natural “barcodes” to help discriminate microbial species from each other, help associate mobile genetic elements to their host-genomes and achieve a more precise shotgun MG analyses.

The ability to link mobile genetic elements to their bacterial hosts potentially allows scientists **to more accurately predict the virulence, antibiotic resistance, and other biologically and clinically critical traits** of individual bacterial species and strains.

SOMATIC VARIANT DETECTION IN COMPLEX CANCER CELL POPULATIONS

Resolve the full spectrum of genetic variation

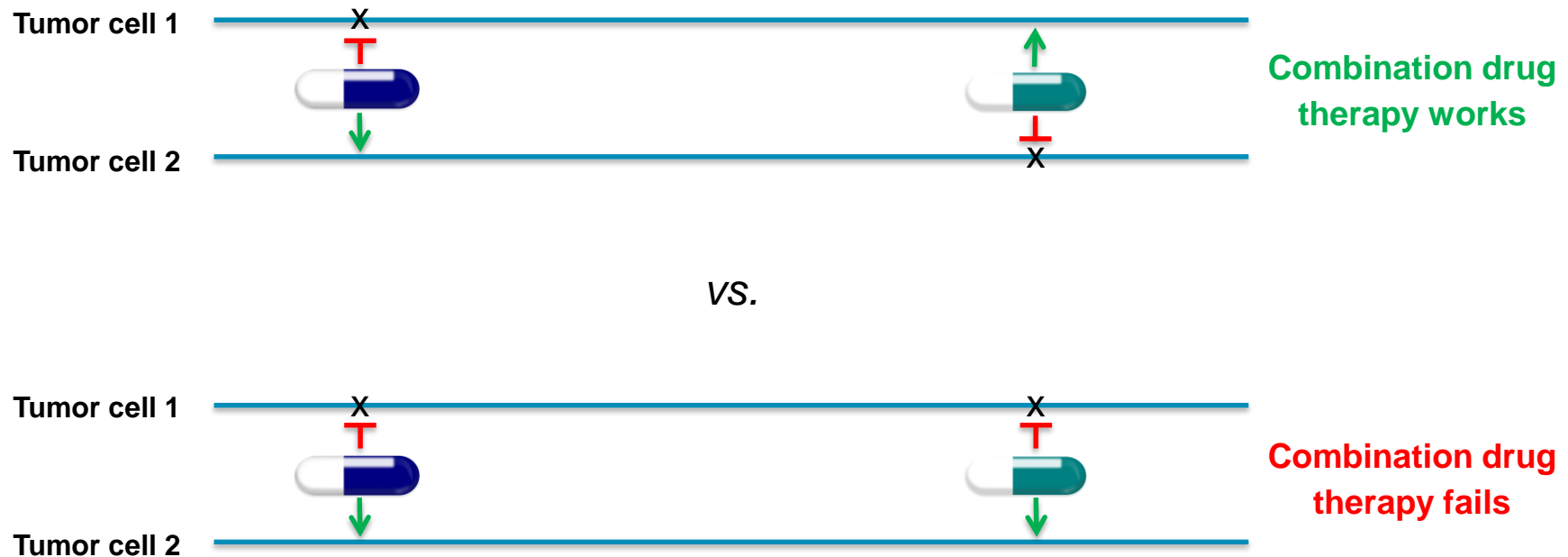
- Tumor samples often contain many unique and evolving genomes, and understanding how complex cancer cell populations adapt in response to treatment is a major focus of cancer research
- SMRT Sequencing delivers a complete, accurate view of the clonal distribution of mutations in genes or genomic regions of interest
- PacBio produces reads long enough to span full-length transcripts and to allow for the immediate detection of compound mutations and splice isoforms. You have the ability to:
 - Detect SNVs occurring at a frequency as low as 1%
 - Differentiate polyclonal from compound mutations
 - View splice variants in gene fusion transcripts
 - Accurately characterize breakpoints in regions of genomic instability



- Identified several mutations and novel transcript isoforms that **standard clinical assays had missed**.
- PacBio sequencing provided clonal distribution frequencies for **compound mutations** and isoforms.

IMPORTANCE OF VARIANT PHASING

Example: Accurate Characterization of Compound Drug Resistance Mutations



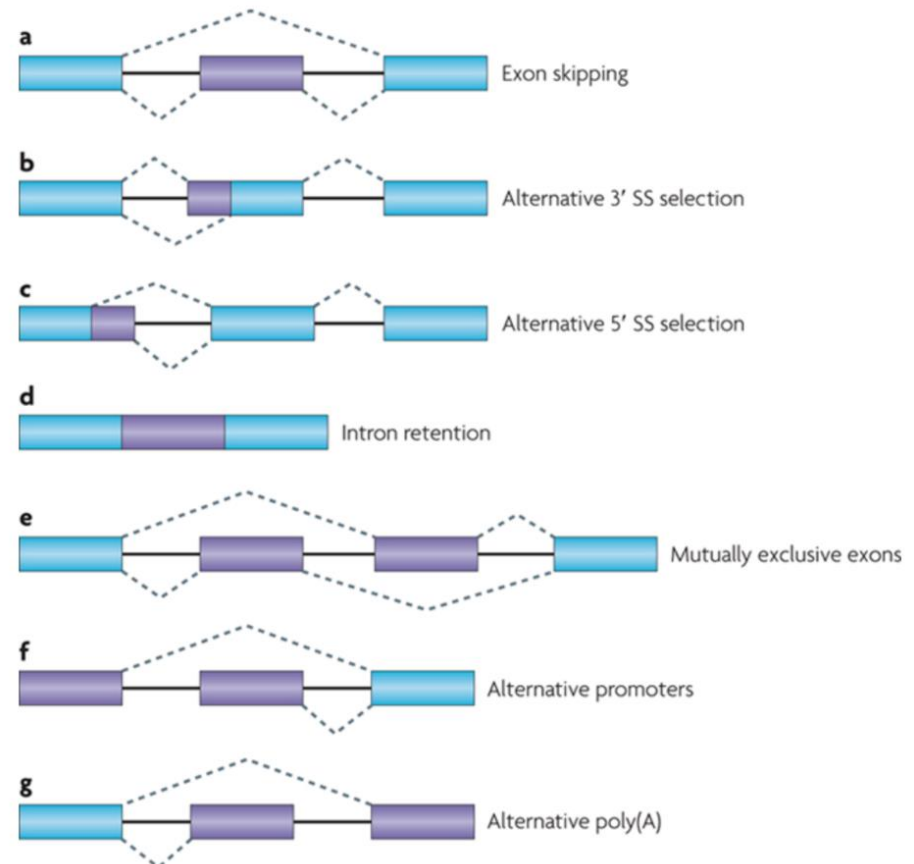


RNA Sequencing (Iso-Seq Method)

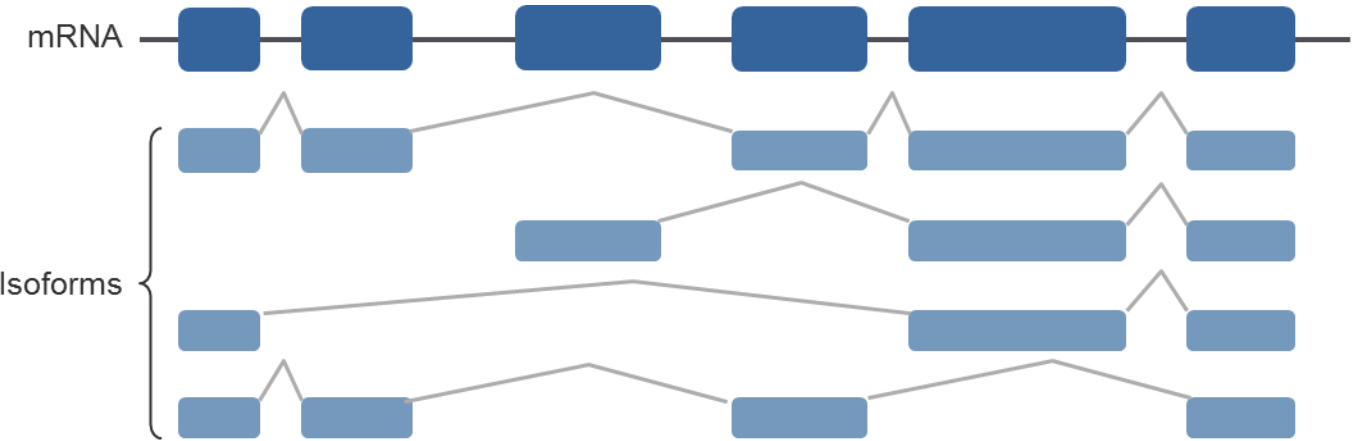
ALTERNATIVE SPLICING OF mRNA

Alternative splicing gives rise to a highly diverse set of proteins from a relatively small number of genes

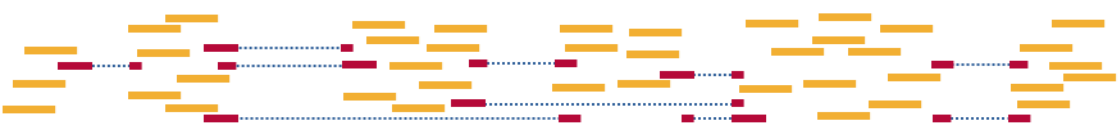
- More than 90% of all genes with multiple exons are alternatively spliced
- Alternative splicing can be divided into different categories:
 - Exon skipping (40%)
 - Alternative 3' splice site selection (18.4%)
 - Alternative 5' splice site selection (7.9%)
 - Intron retention (<5%)
 - Mutually exclusive exons (rare)
 - Alternative promoters (rare)
 - Alternative polyadenylation (rare)



RESOLVING TRANSCRIPTS WITH SHORT READS VS. LONG READS



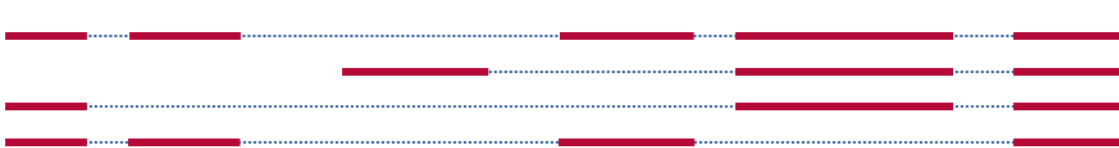
Short-read technologies:



Short reads spanning splice junctions

Insufficient Connectivity
Splice Isoform Uncertainty

PacBio's Iso-Seq™ solution:

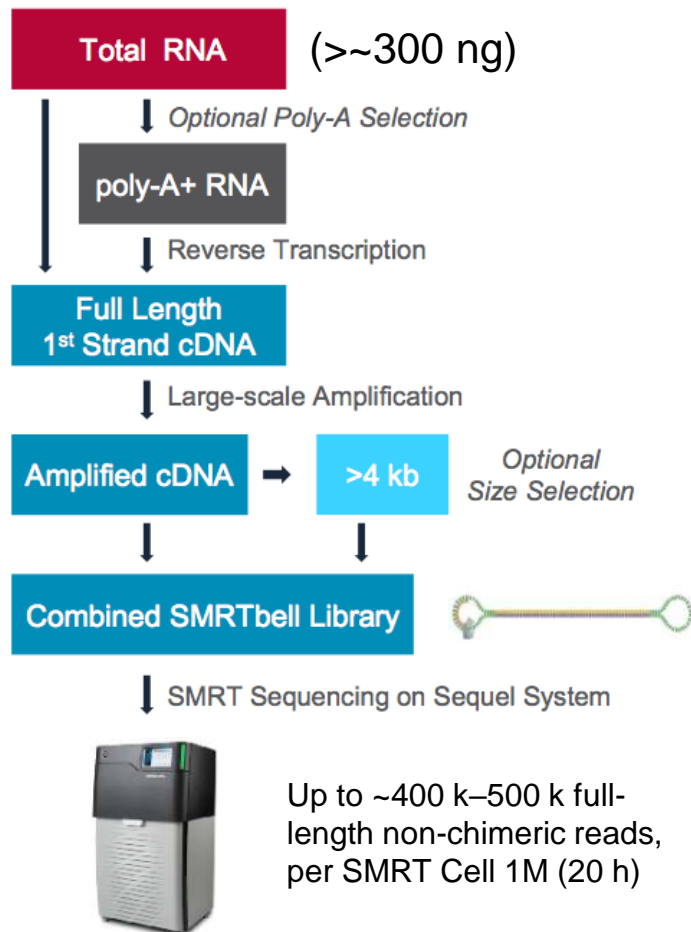


Long reads spanning multiple splice junctions

Full-length cDNA Sequence Reads
Splice Isoform Certainty – No Assembly Required

DISCOVER FULL-LENGTH TRANSCRIPTS WITH THE ISO-SEQ METHOD

Get a complete view of transcript isoform diversity with PacBio long-read sequencing



Iso-Seq Advantages

- *De novo* (reference genome not required)
- No assembly required
- Full-length (5' to 3')
- High accuracy (>99%)

Iso-Seq Applications

- Discover new genes, transcripts and alternative splicing events
- Improve genome annotation to identify gene structure, regulatory elements, and coding regions
- Increase the accuracy of RNA-seq quantification with isoform-level resolution
- Observe allele-specific gene expression

SMRT Sequencing Advantages



Surveying transcript diversity can be done either broadly (**whole transcriptome**) or in a **targeted** fashion

SEQUENCING THE HUMAN CANCER GENOME & TRANSCRIPTOME

Comprehensive SV discovery in the breast cancer cell line SK-BR3

- SK-BR-3 cell line is one of the most important models for HER2+ breast cancers
- SMRT Sequencing of genomic DNA revealed nearly 20,000 structural variants, most of which were missed by short read sequencing
- Full-length transcriptome sequencing using PacBio further revealed several novel gene fusions within nested genomic variants
- Comparison of short and long read technologies for cancer genome analysis revealed a significant gap in our knowledge

Genome Research (2018) 28:1126–1135

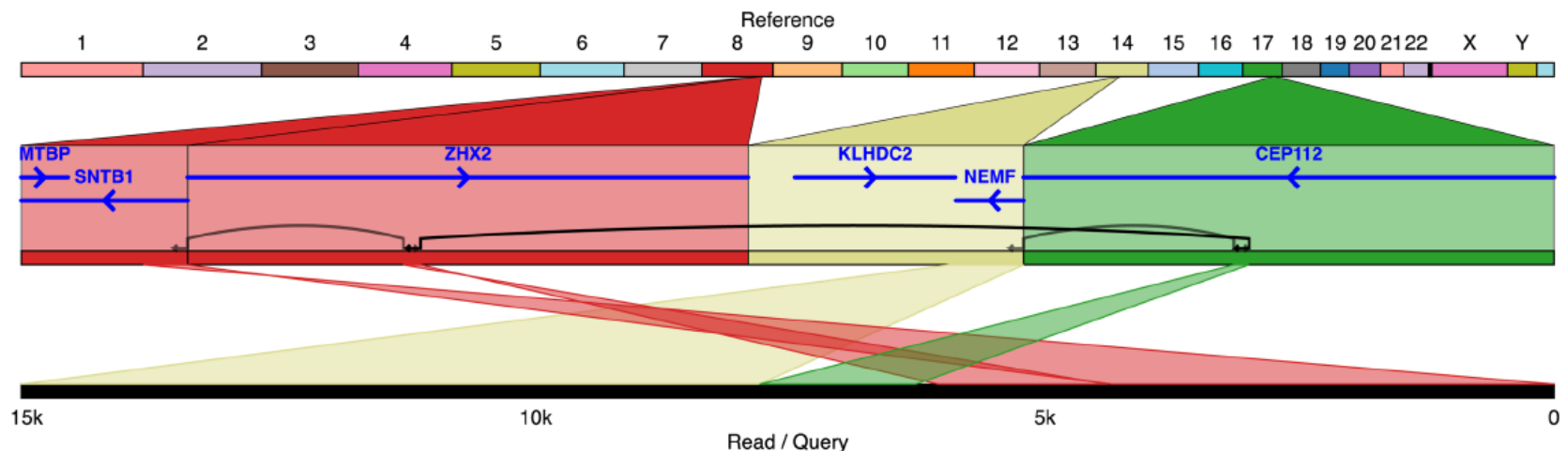


Research

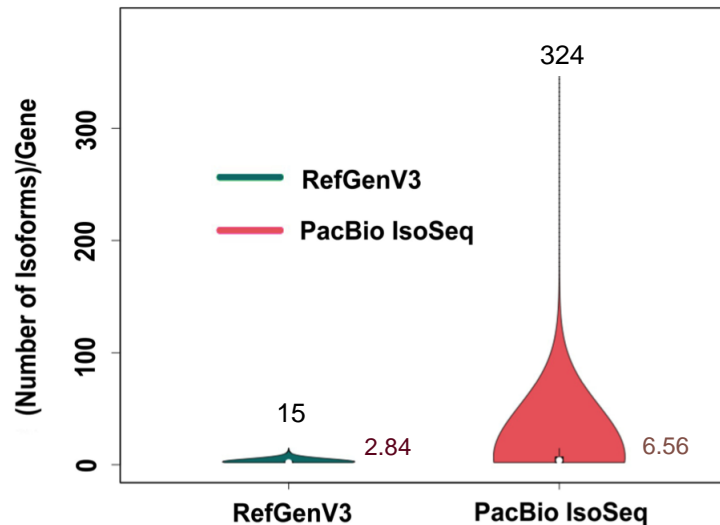
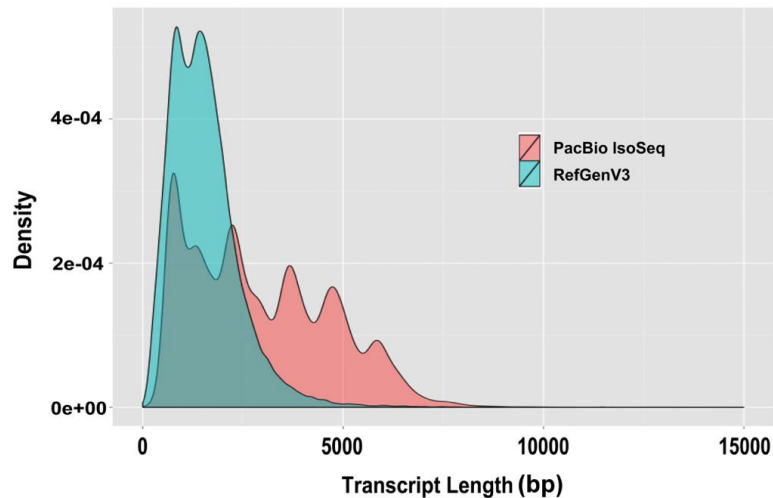
Complex rearrangements and oncogene amplifications revealed by long-read DNA and RNA sequencing of a breast cancer cell line

Maria Nattestad,¹ Sara Goodwin,¹ Karen Ng,² Timour Baslan,³ Fritz J. Sedlazeck,^{4,5} Philipp Rescheneder,⁶ Tyler Garvin,¹ Han Fang,¹ James Gurtowski,¹ Elizabeth Hutton,¹ Elizabeth Tseng,⁷ Chen-Shan Chin,⁷ Timothy Beck,² Yogi Sundaravadanam,² Melissa Kramer,¹ Eric Antoniou,¹ John D. McPherson,⁸ James Hicks,¹ W. Richard McCombie,¹ and Michael C. Schatz^{1,4}

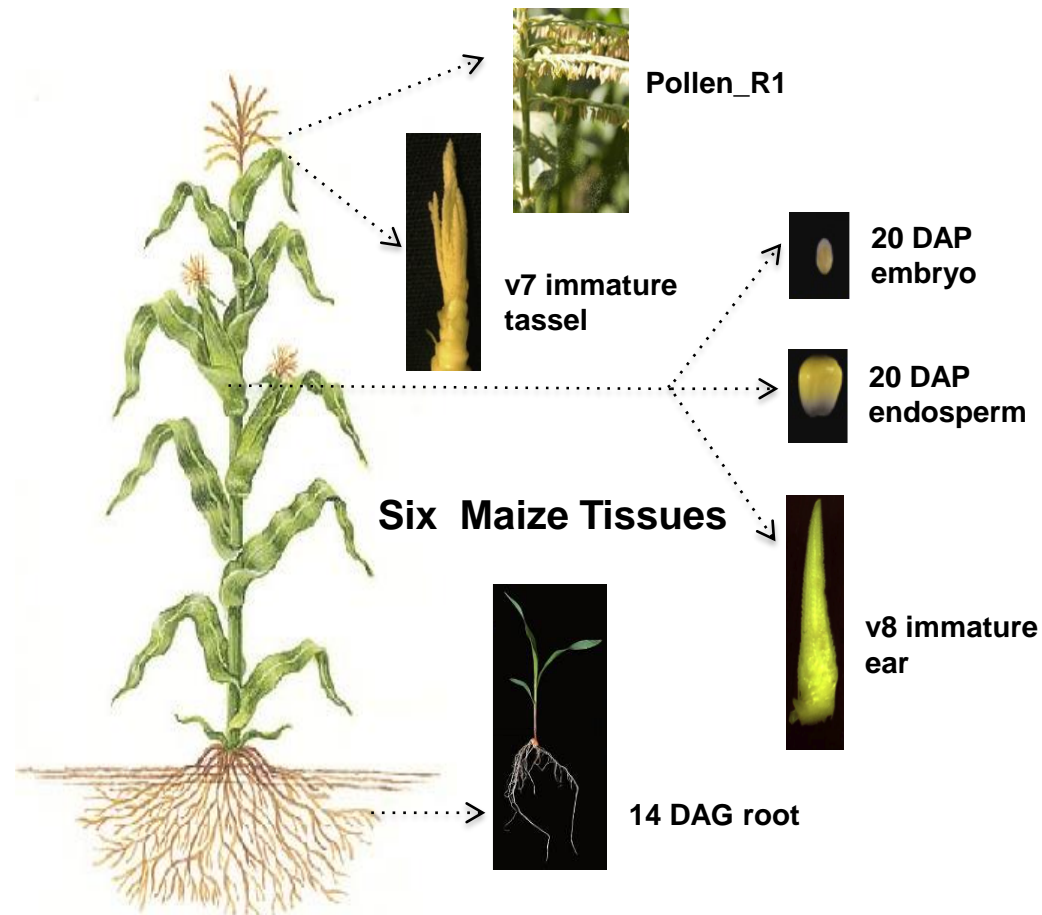
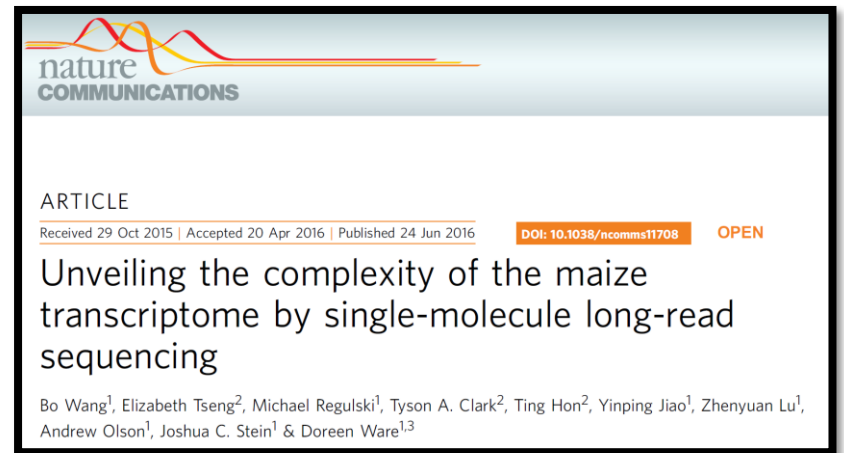
- *SplitThreader* data visualization tool identified **15 high-confidence novel gene fusions** with RNA evidence from the Iso-Seq method and genomic SV evidence from SMRT DNA sequencing
- Discovered a **novel 3-hop gene fusion** between *KLHDC2* and *SNTB1* involving three chromosomes (Chr 8, 14, and 17)



MAIZE B73 ANNOTATION USING MULTIPLEXED ISO-SEQ METHOD

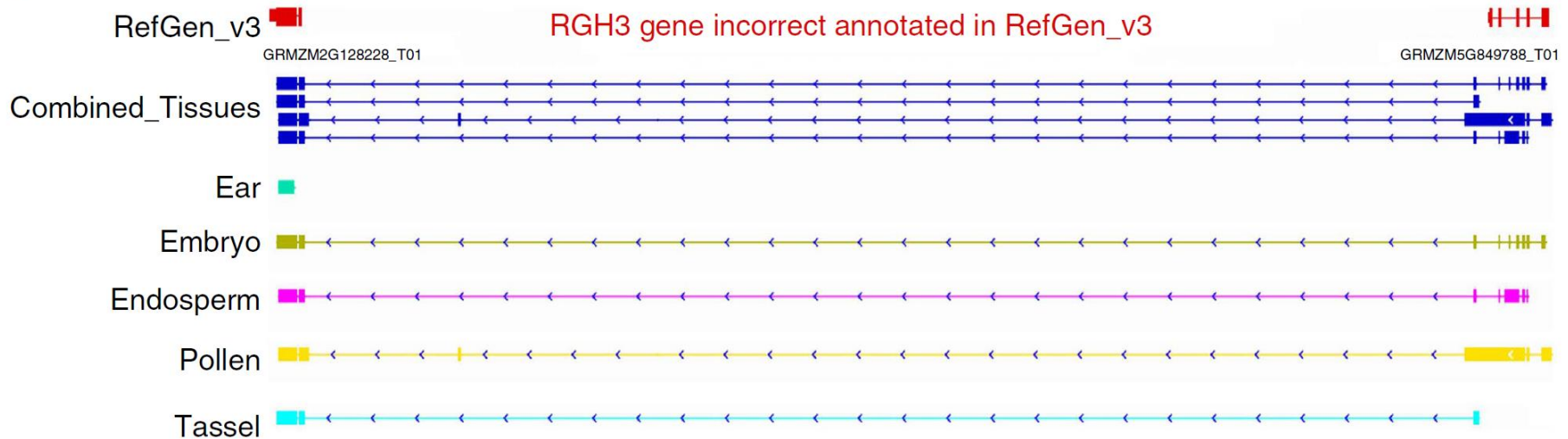


Wang, B. et al. Unveiling the complexity of the maize transcriptome by single-molecule long-read sequencing. Nat Comms 7, 11708 (2016).

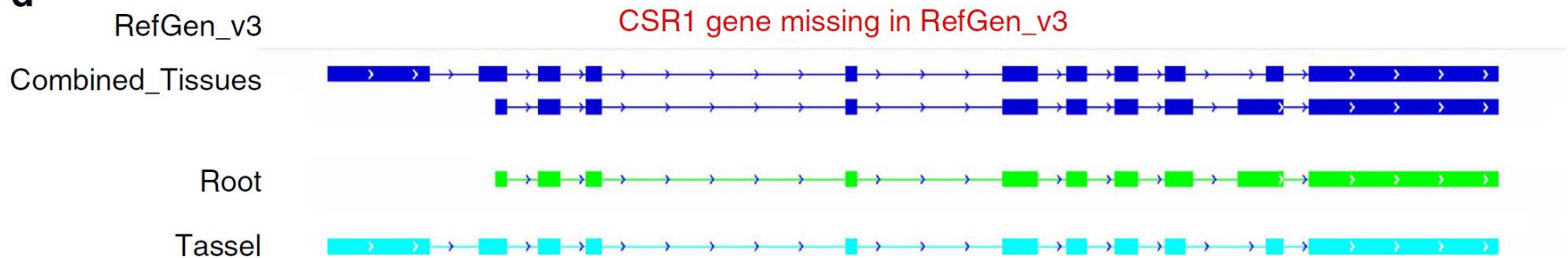


PACBIO ISO-SEQ DATA CORRECTS GENE MODELS FROM THE PREVIOUS MAIZE B73 REFERENCE GENOME V3 ANNOTATION

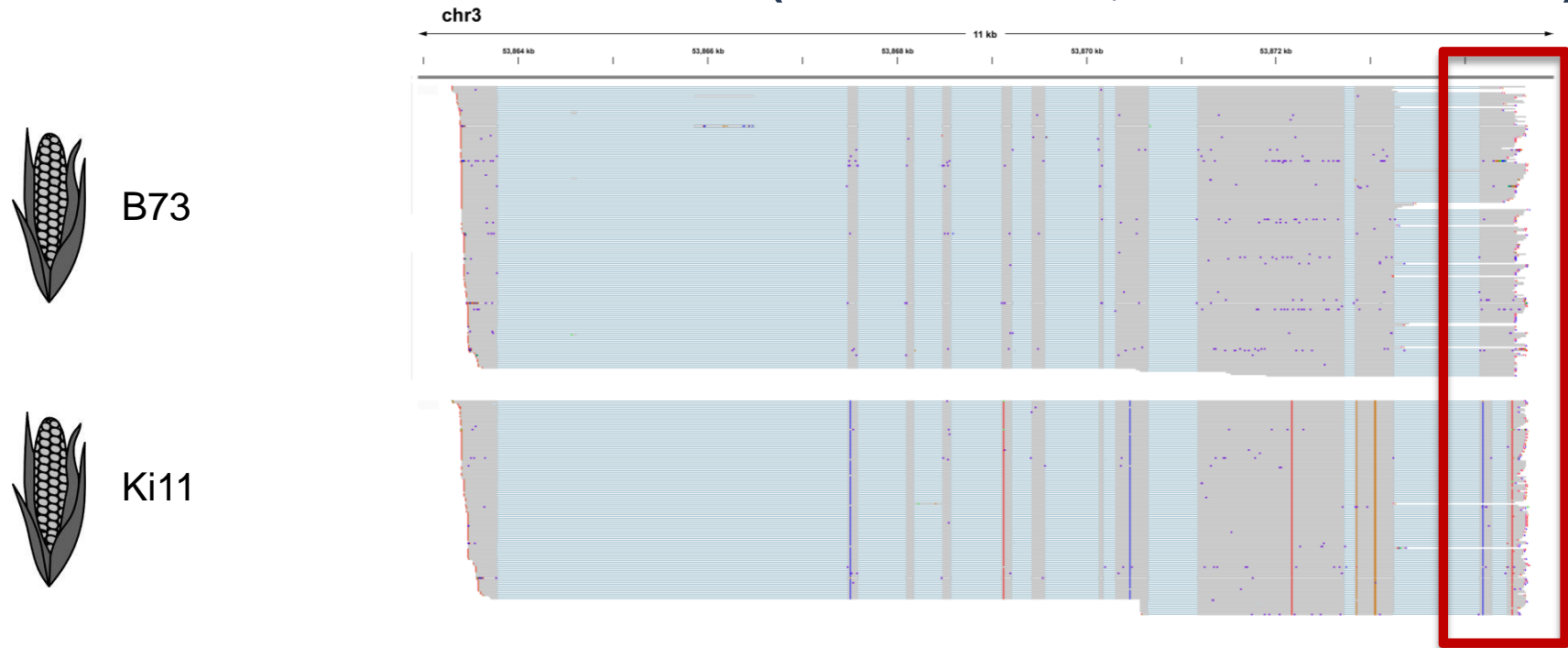
c



d



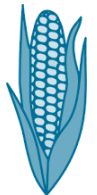
ISO-SEQ ANALYSIS OF ALLELE-SPECIFIC ISOFORM EXPRESSION IN MAIZE F1 HYBRID OFFSPRING (*WANG ET AL., IN PREPARATION*)



male B73

X

female Ki11



male Ki11

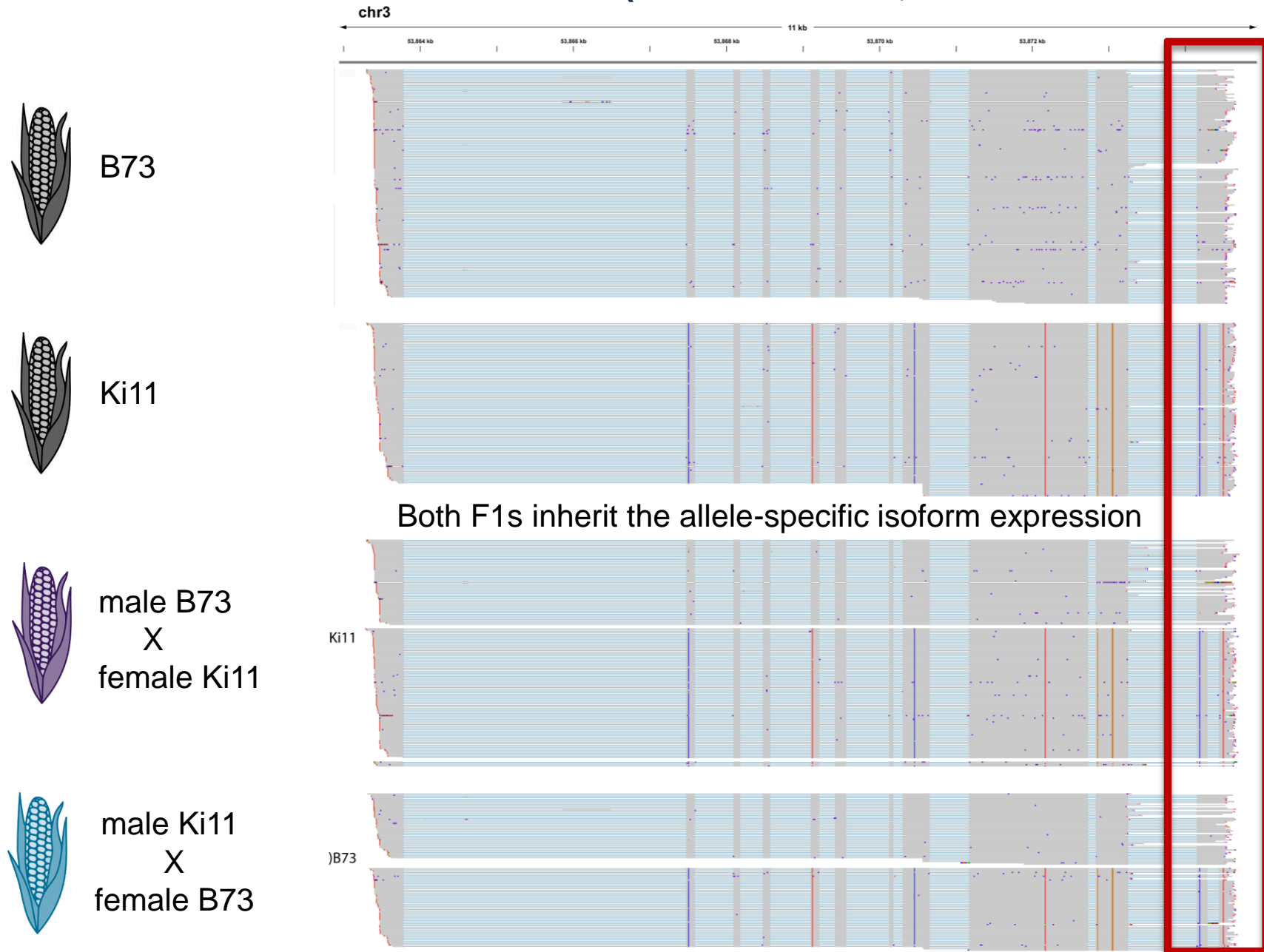
X

female B73

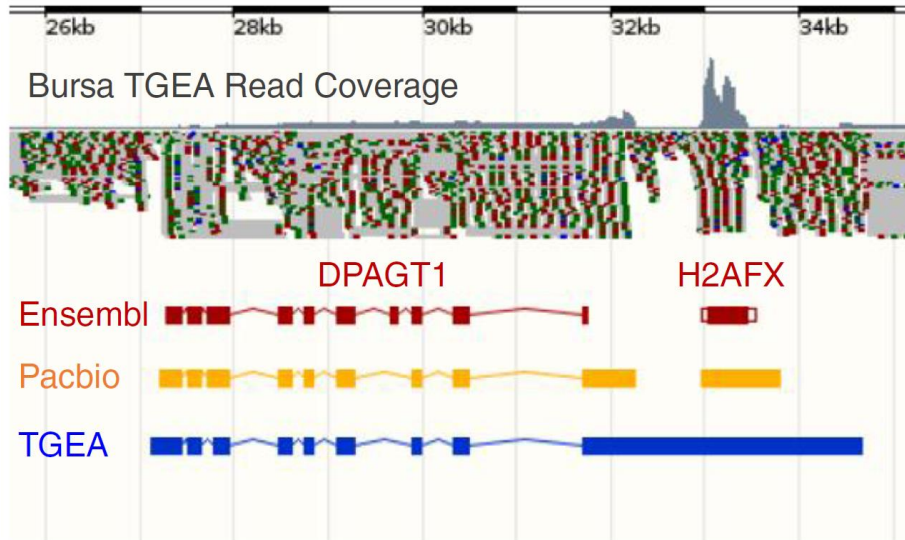
Parent B73 and Ki11 express different isoforms (3' exon difference)

- Two dominant isoforms PB.8517.4 and PB.8517.1
- PB.8517.4 is the canonical isoform and has 11 exons
- PB.8517.1 is a novel isoform with the last exon spliced
- B73 **only** expresses PB.8517.4 (**unspliced** 3' exon)
- Ki11 **only** expresses PB.8517.1 (**spliced** 3' exon)

ISO-SEQ ANALYSIS OF ALLELE-SPECIFIC ISOFORM EXPRESSION IN MAIZE F1 HYBRID OFFSPRING (*WANG ET AL., IN PREPARATION*)



CHICKEN NON-CODING RNA DISCOVERY



Kuo et al. *BMC Genomics* (2017) 18:323
DOI 10.1186/s12864-017-3691-9

BMC Genomics

RESEARCH ARTICLE

Open Access

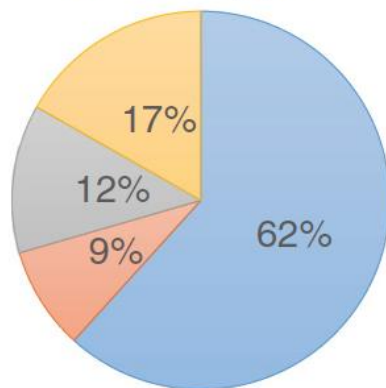
Normalized long read RNA sequencing in chicken reveals transcriptome complexity similar to human

Richard I. Kuo¹, Elizabeth Tseng², Lei Eory¹, Ian R. Paton¹, Alan L. Archibald¹ and David W. Burt^{1,3*}

CrossMark

- Tissue gene expression atlas (TGEA) annotation derived from short-read data **mis-assembled** two genes into one
- PacBio Iso-Seq data unambiguously showed **individual** genes

a Chicken PacBio IncRNA Types Human IncRNA Types Mouse IncRNA Types

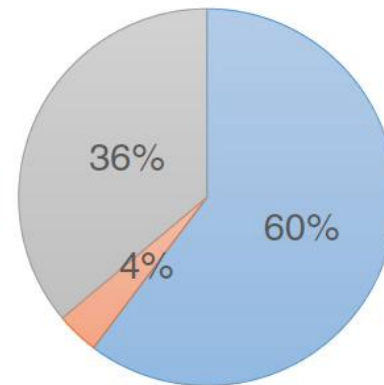
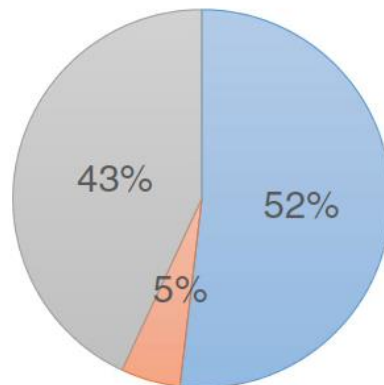


lincRNA

Sense Intronic

Antisense

Sense Exonic



- Normalization increased discovery of low-abundance lncRNAs
- Iso-Seq data showed **chicken lncRNAs are just as diverse and abundant as in human & mouse**



SINGLE CELL ISO-SEQ ANALYSES

Karlsson et al. (2017) BMC Genomics 18: 126

nature|methods

G&T-seq: parallel sequencing of single-cell genomes and transcriptomes

Iain C Macaulay¹, Wilfried Haerty^{2,10}, Parveen Kumar^{3,10}, Yang I Li^{2,9}, Tim Xiaoming Mabel J Teng⁴, Mubeen Goolam⁵, Nathalie Saura Paul Coupland⁷, Lesley M Shirley⁷, Miriam Smit Niels Van der Aa³, Ruby Banerjee⁸, Peter D Ellis Michael A Quail⁷, Harold P Swerdlow^{7,9}, Magdalena Zernicka-Goetz⁵, Frederick J Livesey Chris P Ponting^{1,2,11} & Thierry Voet^{1,3,11}

Macaulay et al. (2015) Nature Methods 12: 519

nature
biotechnology

Letter | Published: 15 October 2018

Single-cell isoform RNA sequencing characterizes isoforms in cerebellar cells

Ishaan Gupta, Paul G Collier, Bettina G. Collier, Ahmed Mahfouz, Anoushka Joglekar, Taylor Floyd, Frank Koopmans, Ben Barres, Augustus Smit, Steven A Sloan, Wenjie Luo, Olivier Fedrigo, M Elizabeth Ross & Hagen U Tilgner

Gupta et al. (2018) Nature Biotechnology 36: 1197

Karlsson and Linnarsson BMC Genomics (2017) 18:126
DOI 10.1186/s12864-017-3528-6

BMC Genomics

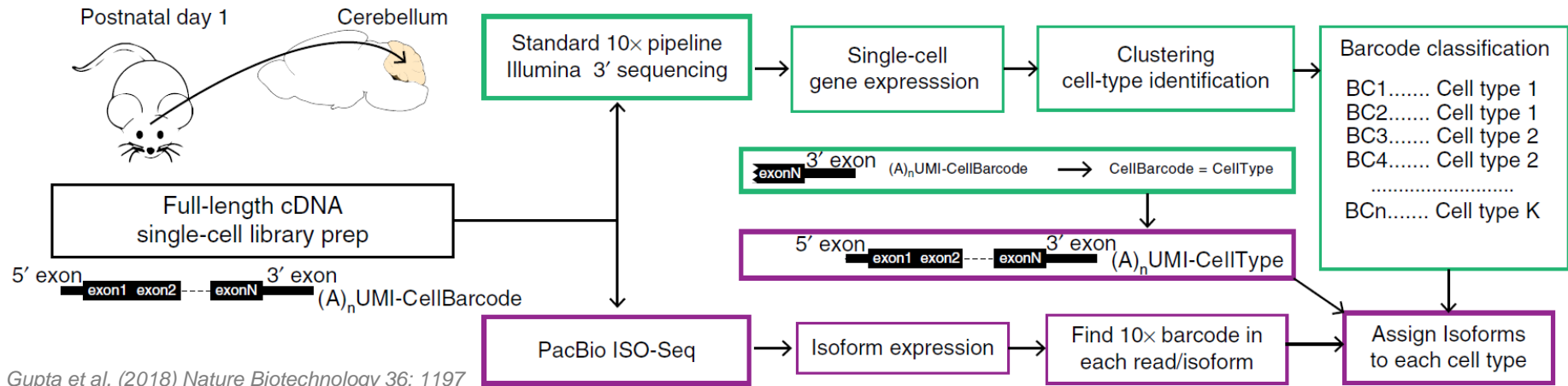
RESEARCH ARTICLE

Open Access

Single-cell mRNA isoform diversity in the mouse brain

Kasper Karlsson¹ and Sten Linnarsson

"We used **ScISO-Seq** to improve genome annotation in mouse Gencode v10 by determining the **cell-type-specific expression of 18,173 known and 16,872 novel isoforms**"



Gupta et al. (2018) Nature Biotechnology 36: 1197



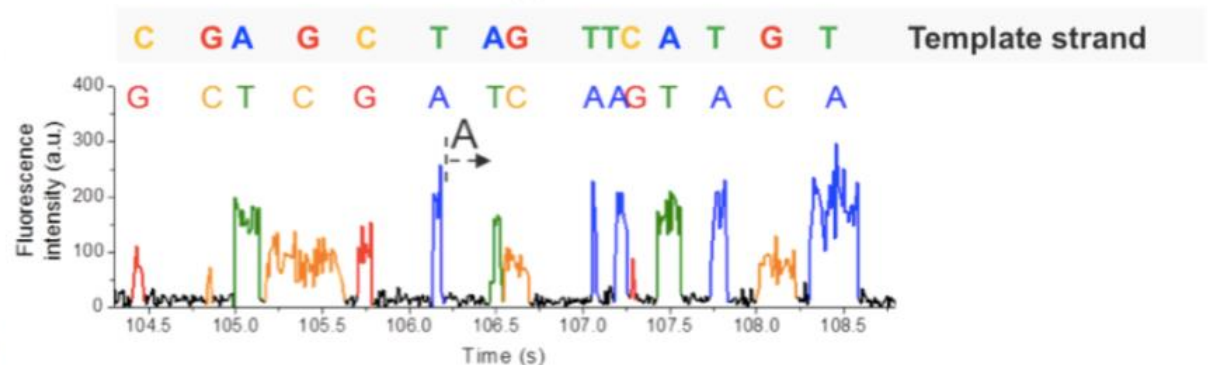
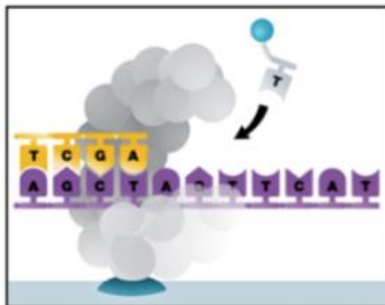
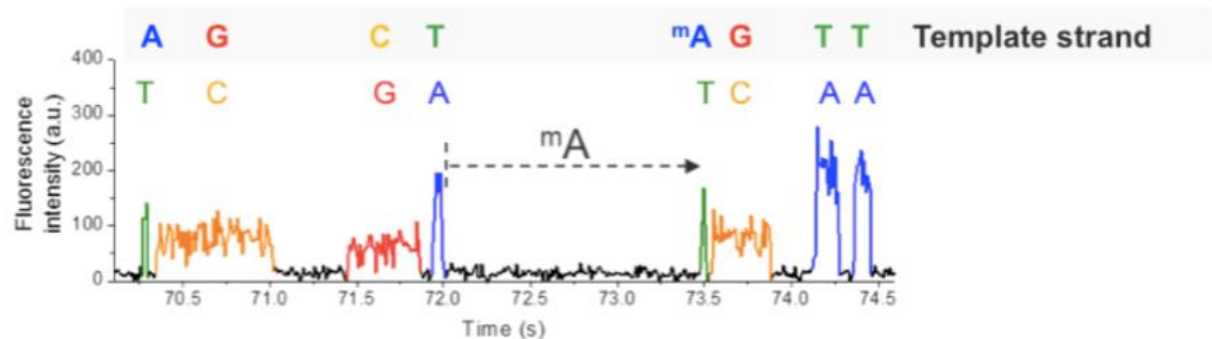
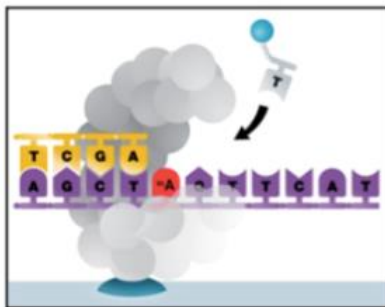
Epigenetics

CHARACTERIZE THE EPIGENETIC LANDSCAPE OF YOUR GENOME

Directly detect epigenetic changes during sequencing to open the door to easier exploration of DNA modification to connect genotype with phenotype.

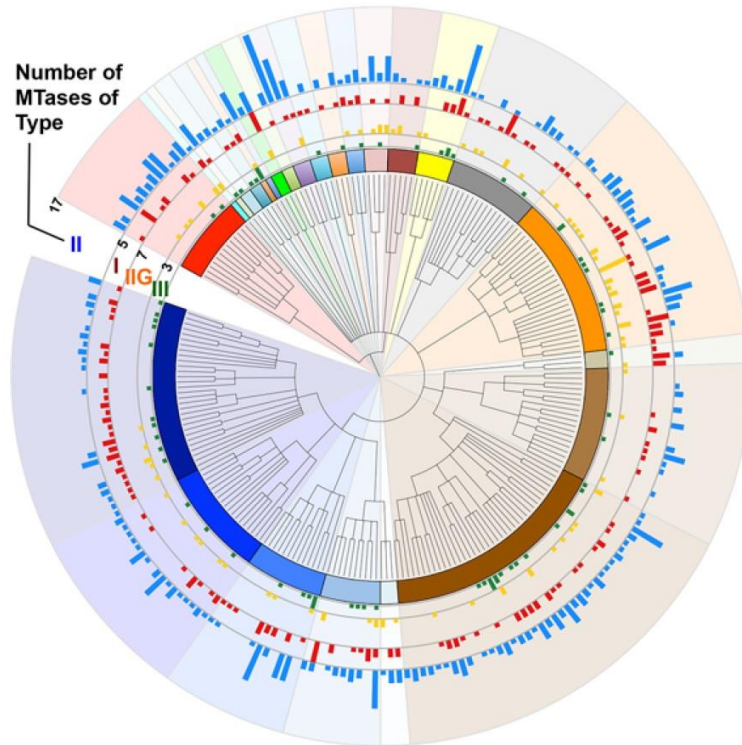
SMRT Sequencing **directly detects epigenetic modifications** by measuring kinetic variation during base incorporation and **eliminates the need for special sample preparation** and additional sequencing.

SMRT Sequencing Advantages

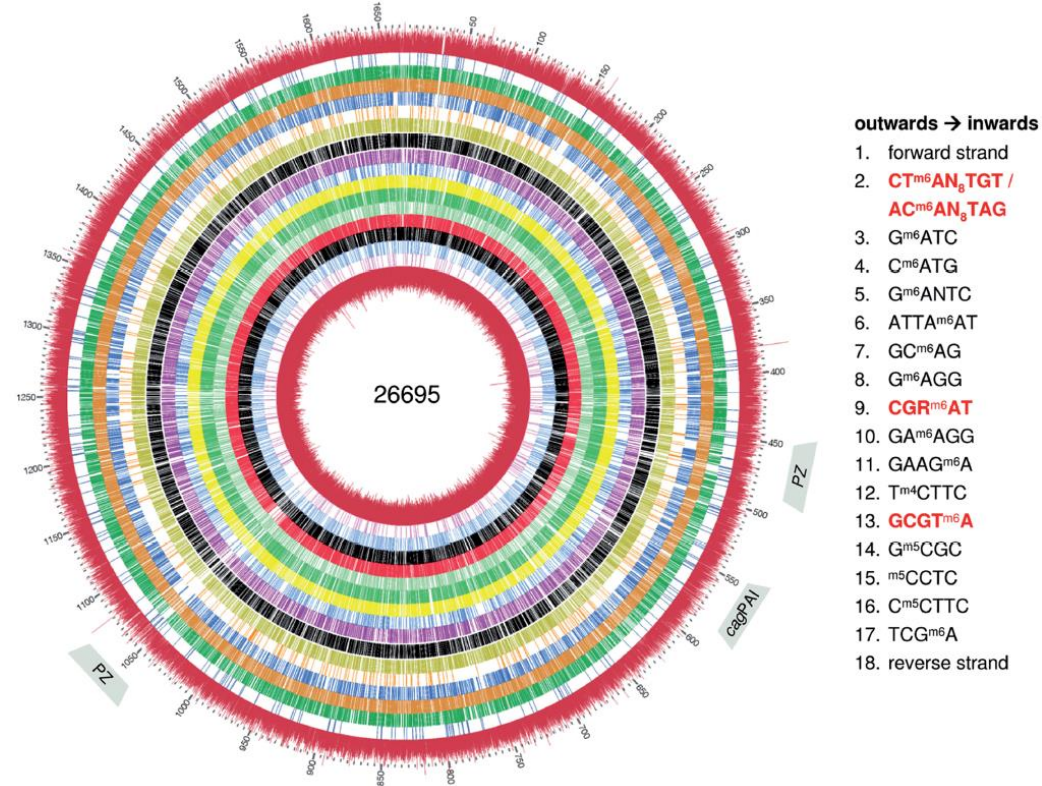


EXPLORE MICROBIAL WHOLE EPIGENOMES – A NEW FRONTIER IN PROKARYOTIC BIOLOGY

- Obtain complete genomes with annotations for epigenetic modification
- Detect genome-wide m6A and m4C R-M system motifs at coverage levels recommended for assembly
- Reveal phase variation of R-M genes that regulate batteries of genes involved in pathogenesis, host adaption, and antibiotic resistance
- Detect strand-specific modification such as hemi-methylation
- Cluster contigs and associated plasmids in metagenomic communities



Phylogenetic tree of 230 sequenced prokaryotic organisms. Outer bars indicate the number and types of active MTases detected per genome.
PLoS Genet 12(2): e1005854

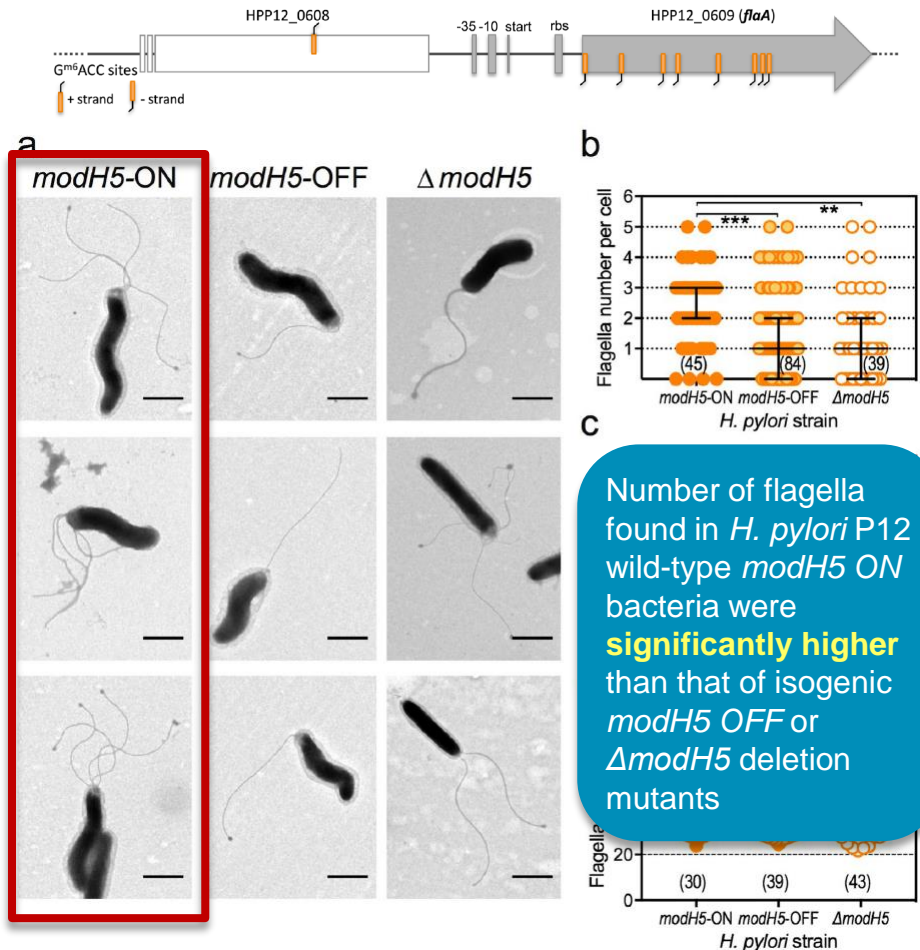
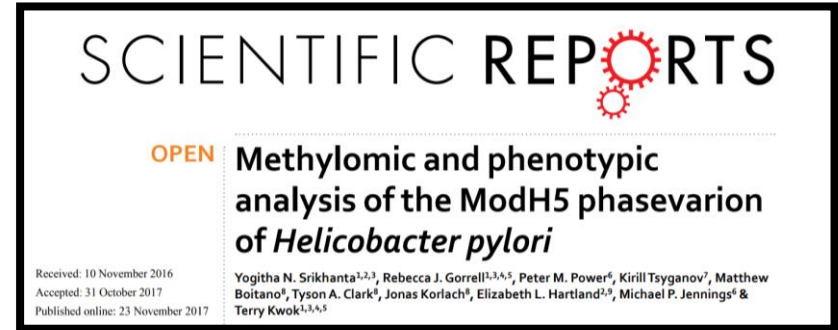
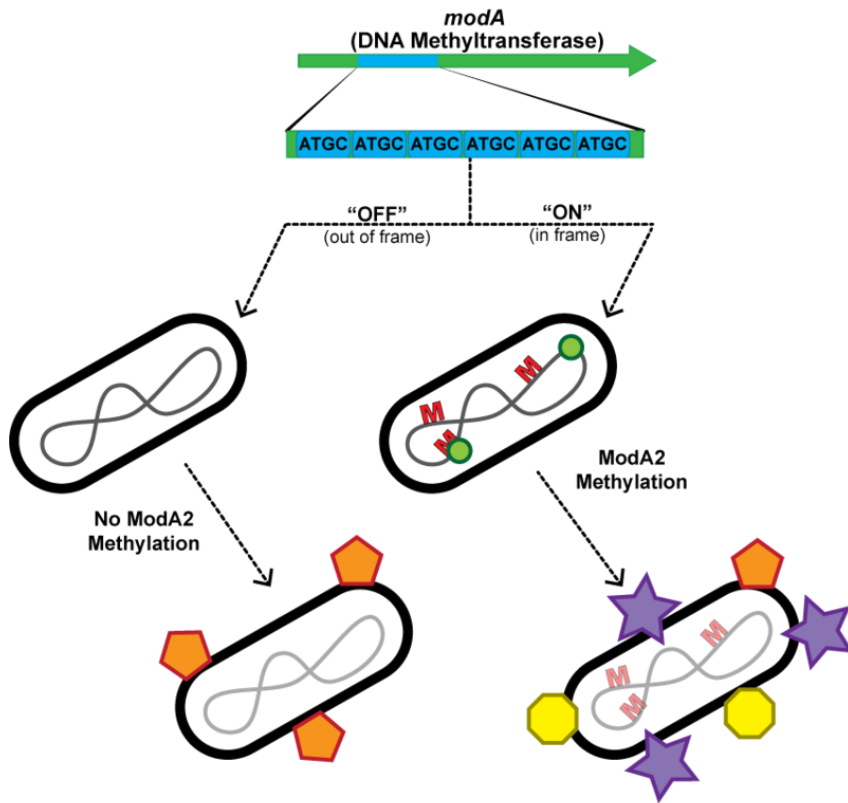


Circos plots displaying the distribution of methylated bases in the genomes of *H. pylori* 26695. The colored tracks in between represent the location of methylation of the different motifs. Novel motifs are highlighted in red in the legend.
Nucleic Acids Research, 2014, Vol. 42, No. 4 2415–2432

PHASE VARIABLE REGULONS OF BACTERIAL PATHOGENS

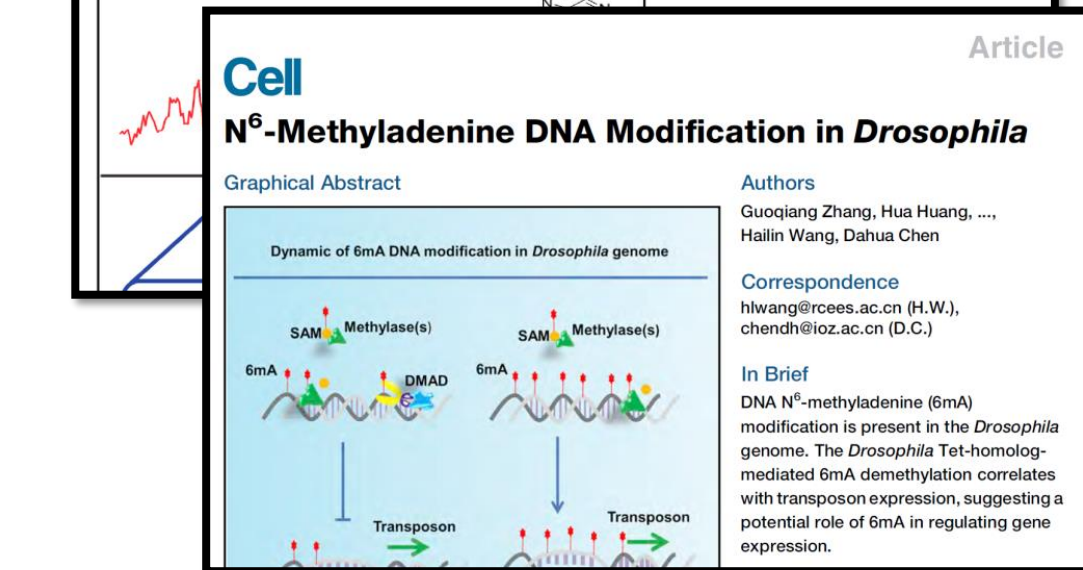
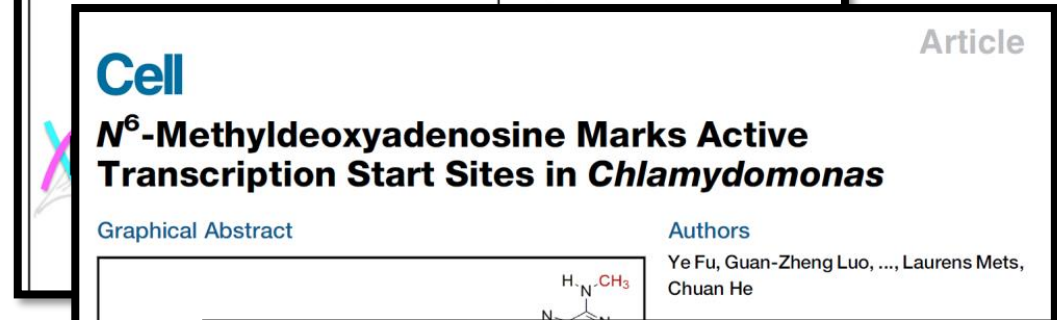
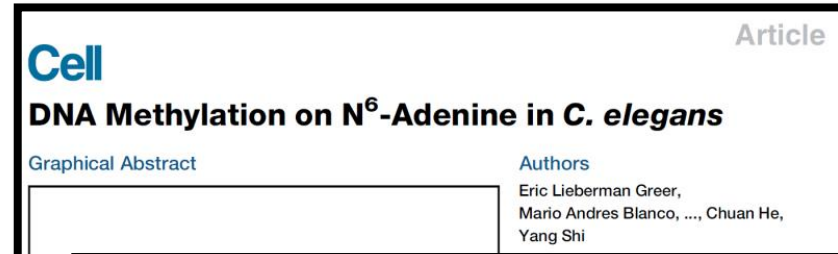
Phasevarions mediate a coordinated change in the expression of multiple bacterial genes or proteins *via* phase variation of a single DNA methyltransferase

Phase variation of methyltransferase expression results in differential methylation throughout the bacterial genome, leading to variable expression of multiple genes through epigenetic mechanisms



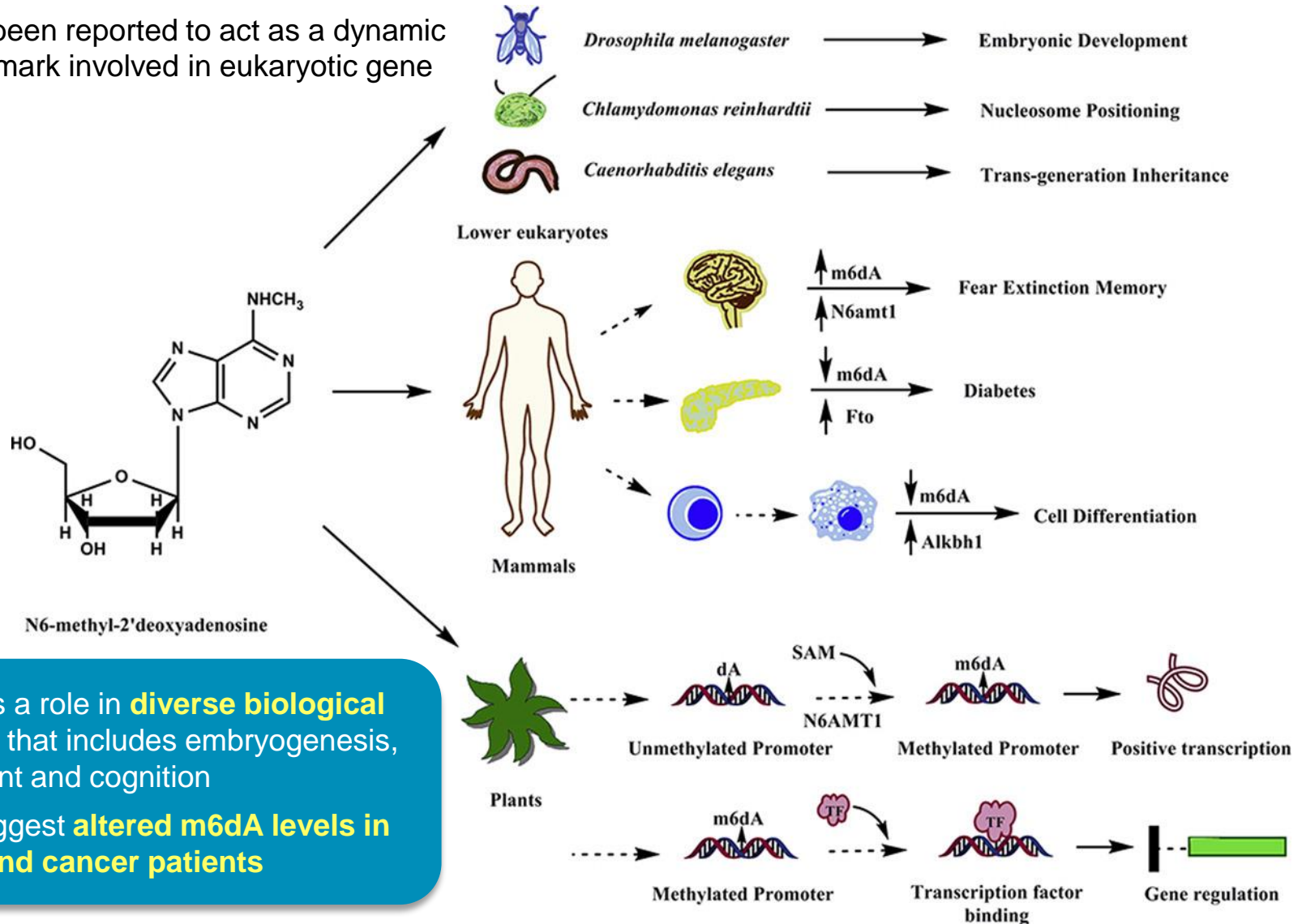
SEE BEYOND A, T, G AND C AND DISCOVER HOW EPIGENETIC MODIFICATIONS PLAY AN IMPORTANT ROLE IN THE BIOLOGY OF EUKARYOTES

- DNA methylation can have an impact on gene expression, imprinting, and X chromosome inactivation.
- Additionally, the deregulation of epigenetic machinery has been implicated in Mendelian disease, human cancer, and drug resistance
- By pairing high throughput with long reads, PacBio offers scalable solutions for assessing CpG methylation in eukaryotic genomes.
- For targeted applications, SMRT bisulfite sequencing marries bisulfite samples with highly multiplexed, quantitative long-read sequencing, accurately measuring CpG methylation across kilobase-sized regions
- N6-methyl-2'-deoxyadenosine (m6dA) has recently been identified in several phylogenetically distinct eukaryotes.



N6-METHYL-2'-DEOXYADENOSINE (M6DA) HAS RECENTLY BEEN IDENTIFIED IN SEVERAL PHYLOGENETICALLY DISTINCT EUKARYOTES

- m6dA has been reported to act as a dynamic epigenetic mark involved in eukaryotic gene regulation



- m6dA plays a role in **diverse biological processes** that includes embryogenesis, development and cognition
- Studies suggest **altered m6dA levels in diabetes and cancer patients**



PACIFIC
BIOSCIENCES®



Summary

SEQUEL SYSTEM: A FOUNDATION FOR DISCOVERY USING SINGLE-MOLECULE, REAL-TIME (SMRT) DNA SEQUENCING

Accelerate your research with the most comprehensive view of genomes, transcriptomes, and epigenomes. Reduce project costs and timelines as you create highest-quality whole genome assemblies and explore the full size-spectrum of genetic variation.

- ❖ Create high-quality whole genome *de novo* assemblies of organisms
- ❖ Cost-effectively survey large population cohort studies for structural variants at low-fold coverage
- ❖ Target hard-to-sequence regions not easily accessible by other technologies
- ❖ Detect genomic variation in complex population mixtures with highly accurate long reads
- ❖ Sequence full-length transcriptomes or targeted transcripts with no assembly required
- ❖ Detect epigenetic modifications without using complicated sample preparation techniques



Win *Free* PacBio sequencing! Local SMRT Grant Now Open

PacBio long-read sequencing services are now offered by the **McMaster Genomics Facility**
fmf@mcmaster.ca

THE LEADER IN LONG-READ SEQUENCING



Single Molecule, Real-Time (SMRT®) Sequencing delivers long read lengths with the highest consensus accuracy and uniform coverage, allowing you to go beyond fragmented draft genomes and transcriptome reconstruction using isoform-inference algorithms. How could your research benefit from long-reads? Let us know and you could win SMRT Sequencing at the [McMaster Genomics Facility](#).

***This Local SMRT Grant is
brought to you by***



Sonny Mark, Ph.D.
Staff Scientist
Pacific Biosciences
smark@pacb.com

Christine Mader
Manager
McMaster Genomics Facility
kingc@mcmaster.ca

- ❖ Submit your 250-word proposal at www.pacb.com/localsmrtgrant by **May 15th, 2019**
- ❖ Winning project receives up to C\$5000 in sequencing services (bioinformatics support not included)

PacBio Office Hours: 1:00 PM – 5:00 PM
HSC 3N50

McMaster University Health Sciences Centre

Online Scheduler: <https://PacBio.as.me/>

Select any open 30-min timeslot on Apr. 11th or email Sonny Mark, Ph.D., PacBio Scientist at smark@pacb.com

Light refreshments will be available during the sessions



www.pacb.com

For Research Use Only. Not for use in diagnostic procedures. © Copyright 2019 by Pacific Biosciences of California, Inc. All rights reserved. Pacific Biosciences, the Pacific Biosciences logo, PacBio, SMRT, SMRTbell, Iso-Seq, and Sequel are trademarks of Pacific Biosciences. BluePippin and SageELF are trademarks of Sage Science. NGS-go and NGSengine are trademarks of GenDx. FEMTO Pulse and Fragment Analyzer are trademarks of Advanced Analytical Technologies.

All other trademarks are the sole property of their respective owners.



PACIFIC
BIOSCIENCES®



Where to Find More Information

SMRT SEQUENCING RESOURCES

Explore our collection of technical resources and learn how scientists use SMRT Sequencing to advance their research

<https://www.pacb.com/smrt-science/smrt-resources/>



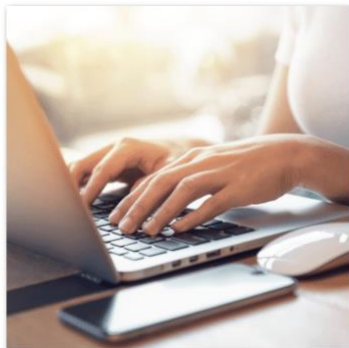
Scientific Publications



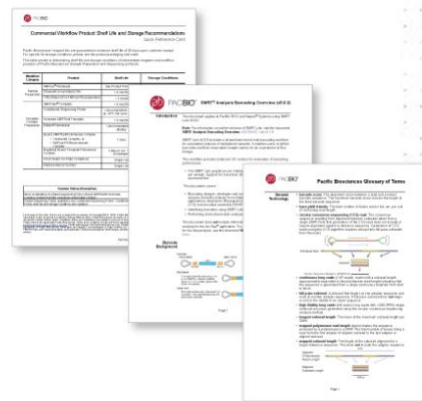
PacBio Literature



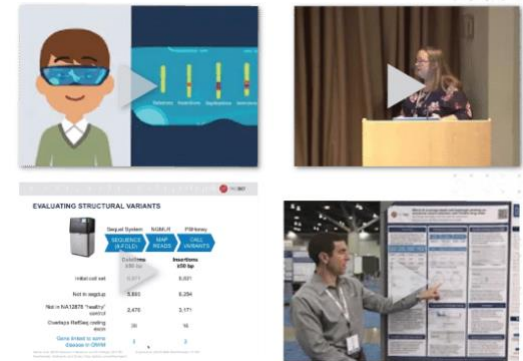
Posters



BLOG



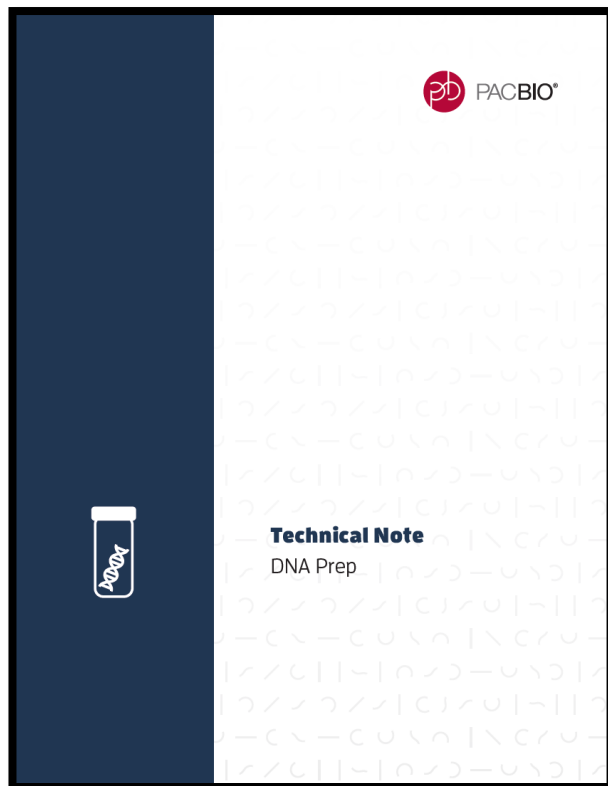
Documentation



Video Gallery

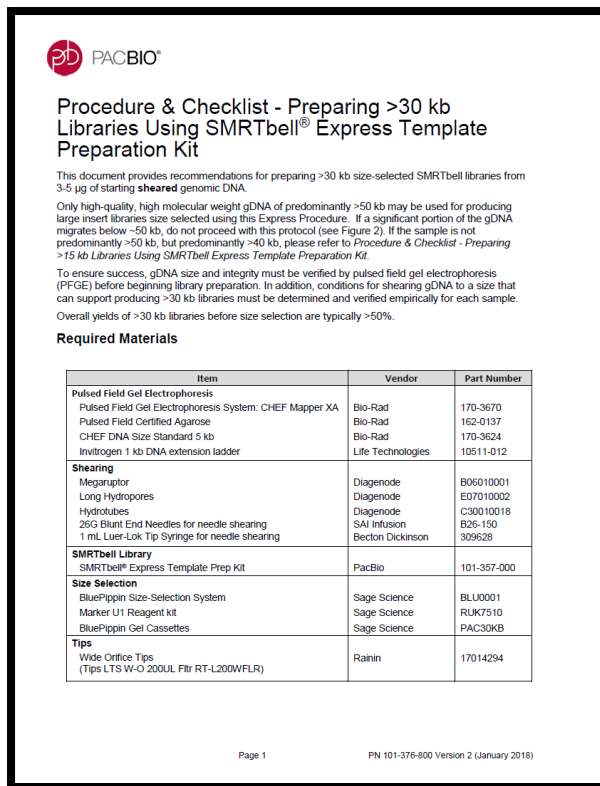
DNA SAMPLE QC AND LIBRARY PREPARATION PROTOCOL DOCUMENTATION

DNA Sample and SMRTbell Library QC Technical Notes



Technical Note DNA Prep

SMRTbell Library Prep Procedures & Checklists



Procedure & Checklist - Preparing >30 kb Libraries Using SMRTbell® Express Template Preparation Kit

This document provides recommendations for preparing >30 kb size-selected SMRTbell libraries from 3-5 µg of starting **sheared** genomic DNA.

Only high-quality, high molecular weight gDNA of predominantly >50 kb may be used for producing large insert libraries size selected using this Express Procedure. If a significant portion of the gDNA migrates below ~50 kb, do not proceed with this protocol (see Figure 2). If the sample is not predominantly >50 kb, but predominantly >40 kb, please refer to Procedure & Checklist - Preparing >15 kb Libraries Using SMRTbell Express Template Preparation Kit.

To ensure success, gDNA size and integrity must be verified by pulsed field gel electrophoresis (PFGE) before beginning library preparation. In addition, conditions for shearing gDNA to a size that can support producing >30 kb libraries must be determined and verified empirically for each sample. Overall yields of >30 kb libraries before size selection are typically >50%.

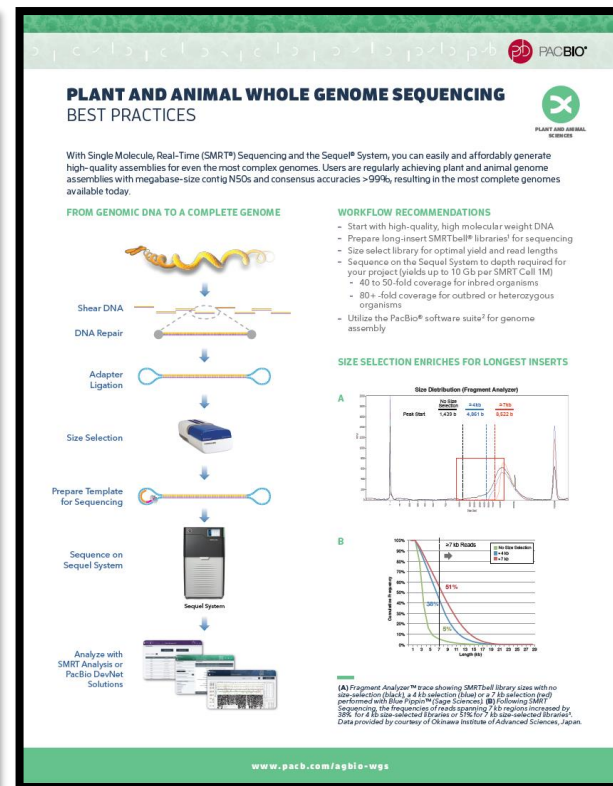
Required Materials

| Item | Vendor | Part Number |
|---|-------------------|-------------|
| Pulsed Field Gel Electrophoresis | | |
| Pulsed Field Gel Electrophoresis System: CHEF Mapper XA | Bio-Rad | 170-3670 |
| Pulsed Field Certified Agarose | Bio-Rad | 162-0137 |
| CHEF DNA Size Standard 5 kb | Bio-Rad | 170-3624 |
| Invitrogen 1 kb DNA extension ladder | Life Technologies | 10511-012 |
| Shearing | | |
| Megaruptor | Diagenode | B06010001 |
| Long Hydrotubes | Diagenode | E07010002 |
| Hydrotubes | Diagenode | C30010018 |
| 26G Blunt End Needles for needle shearing | SAI Infusion | B26-150 |
| 1 mL Luer-Lok Tip Syringe for needle shearing | Bedon Dickinson | 309628 |
| SMRTbell Library | | |
| SMRTbell Express Template Prep Kit | PacBio | 101-357-000 |
| Size Selection | | |
| BluePippin Size-Selection System | Sage Science | BLU0001 |
| Marker U1 Reagent kit | Sage Science | RUUK7510 |
| BluePippin Gel Cassettes | Sage Science | PAC30KB |
| Tips | | |
| Wide Office Tips | Rainin | 17014294 |
| (Tips LTS W-O 200UL FTR RT-L200WFLR) | | |

Page 1

PN 101-376-800 Version 2 (January 2018)

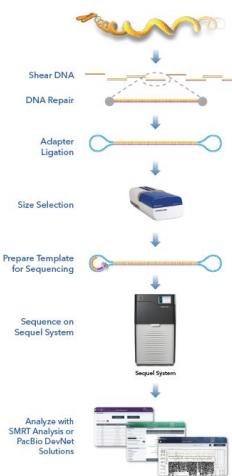
Best Practices Application Briefs



PLANT AND ANIMAL WHOLE GENOME SEQUENCING BEST PRACTICES

With Single Molecule, Real-Time (SMRT®) Sequencing and the Sequel® System, you can easily and affordably generate high-quality assemblies for even the most complex genomes. Users are regularly achieving plant and animal genome assemblies with megabase-size contig N50s and consensus accuracies >99%, resulting in the most complete genomes available today.

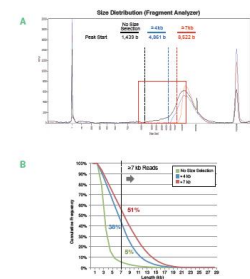
FROM GENOMIC DNA TO A COMPLETE GENOME



WORKFLOW RECOMMENDATIONS

- Start with high-quality, high molecular weight DNA
- Prepare long-insert SMRTbell® libraries for sequencing
- Size select library for optimal yield and read lengths
- Sequence on the Sequel System to depth required for your project (yields up to 10 Gb per SMRT Cell 1M)
- 40 to 50-fold coverage for inbred organisms
- 80+ fold coverage for outbred or heterozygous organisms
- Utilize the PacBio® software suite for genome assembly

SIZE SELECTION ENRICHES FOR LONGEST INSERTS



(A) Fragment Analysis™ trace showing SMRTbell library sizes with no size selection (black) or 4 kb size selection (blue) or 4 kb size selection performed with BluePippin™ (Sage Science). (B) Following SMRT Sequencing, the frequency of reads reporting 1 kb regions increased by 30% for a 4 kb size-selected library or 17% for a 7 kb size-selected library. Data provided by courtesy of Okazaki Institute of Advanced Sciences, Japan.

www.pacb.com/agbio-wgs

<https://www.pacb.com/smrt-science/smrt-resources/pacbio-literature/>

<https://www.pacb.com/support/documentation/>

<https://www.pacb.com/smrt-science/smrt-resources/pacbio-literature/>

<https://www.pacb.com/support/documentation/> and <https://www.pacb.com/smrt-science/smrt-resources/pacbio-literature/>

Genome Assembly Online Project Builder

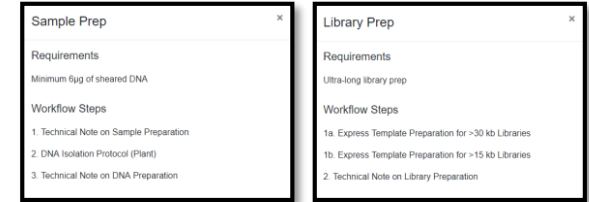
- <https://www.pacb.com/genome-project-builder/>
- Project builder tool for planning *de novo* assembly projects
- Receive tailored recommendations for each step in the workflow from sample prep to sequencing / assembly and annotation.



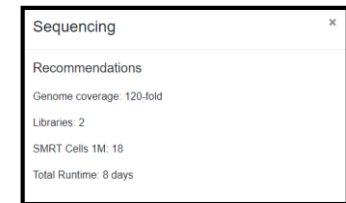
BUILD YOUR PROJECT ➔



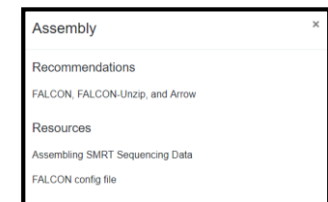
Sample, DNA & Library Prep



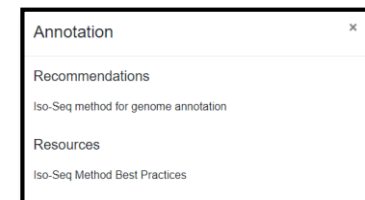
Sequencing



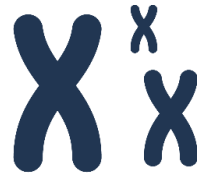
Genome Assembly & Polishing



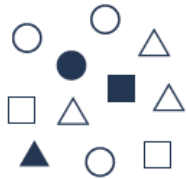
Genome Annotation



THE MOST COMPREHENSIVE VIEW OF GENOMES, TRANSCRIPTOMES, AND EPIGENOMES



**WHOLE GENOME
SEQUENCING**



**COMPLEX
POPULATIONS**



PACBIO®



RNA SEQ



TARGETED SEQ

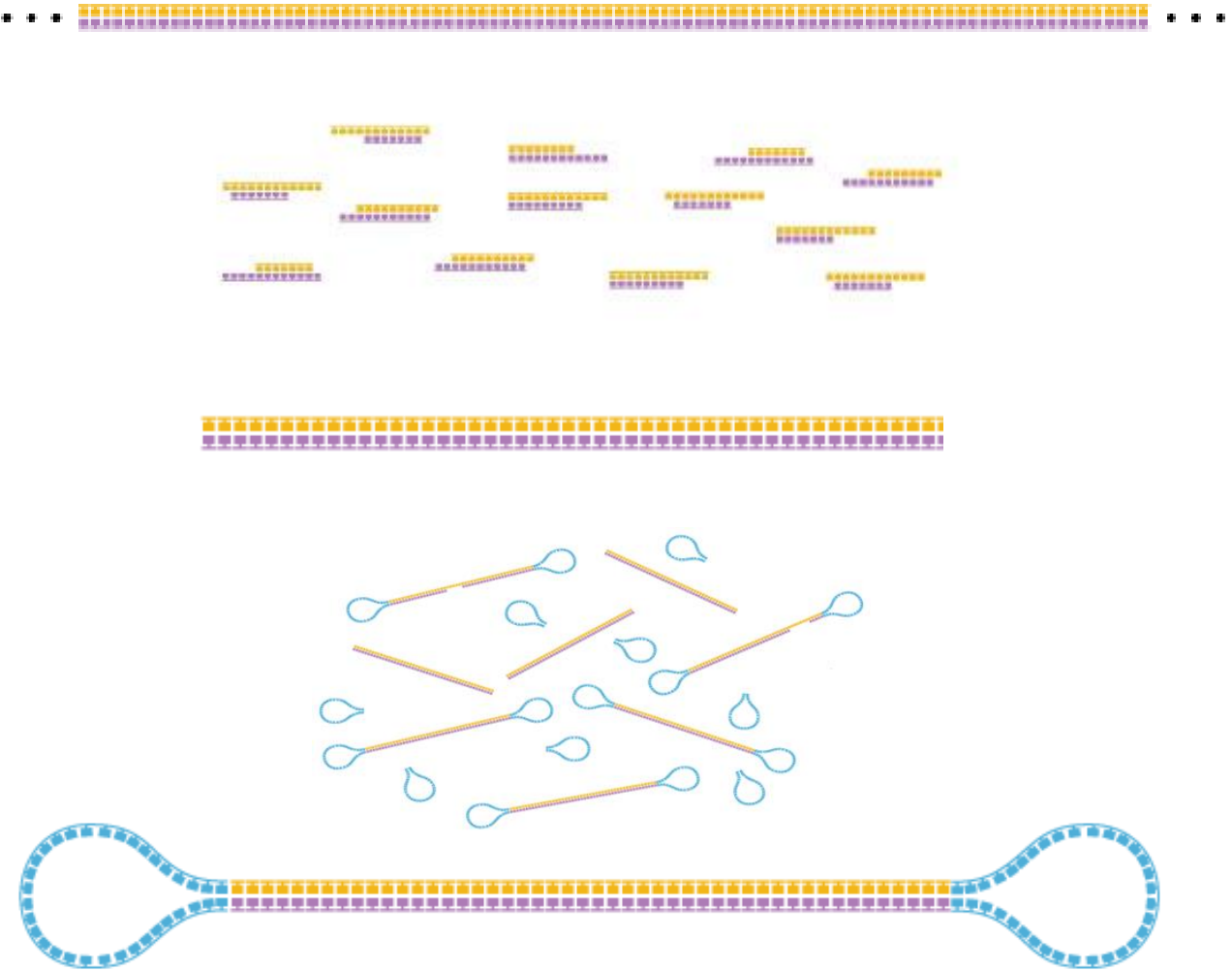


EPIGENETICS

DNA SMRTBELL LIBRARY TEMPLATE PREPARATION OVERVIEW



| Insert Size (bp) | Typical Input DNA per Prep (ng) |
|-------------------------|---------------------------------|
| 250 – 500 | 250 |
| 1,000 – 2000 | 500 |
| 5,000 – 10,000 | 1,000 |
| >15,000 (Size-selected) | 2,000 - 5,000 |



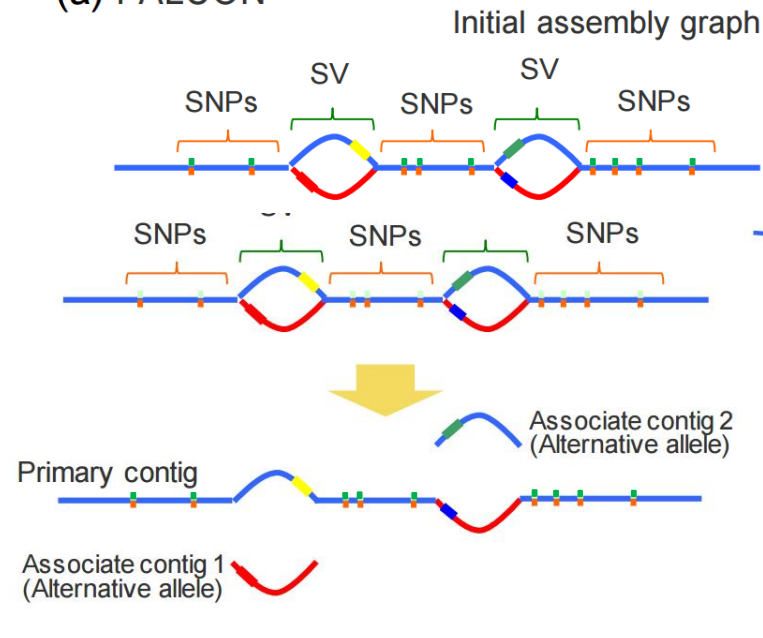
INTERACTIVE MAP OF HUMAN GENOME ASSEMBLIES

>40 publicly available assemblies using PacBio sequencing data

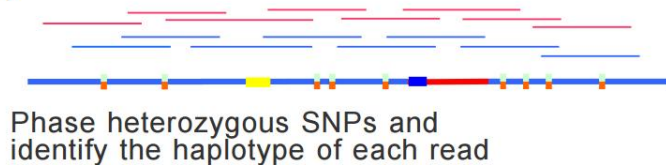


DIPLOID ASSEMBLY WITH FALCON-UNZIP

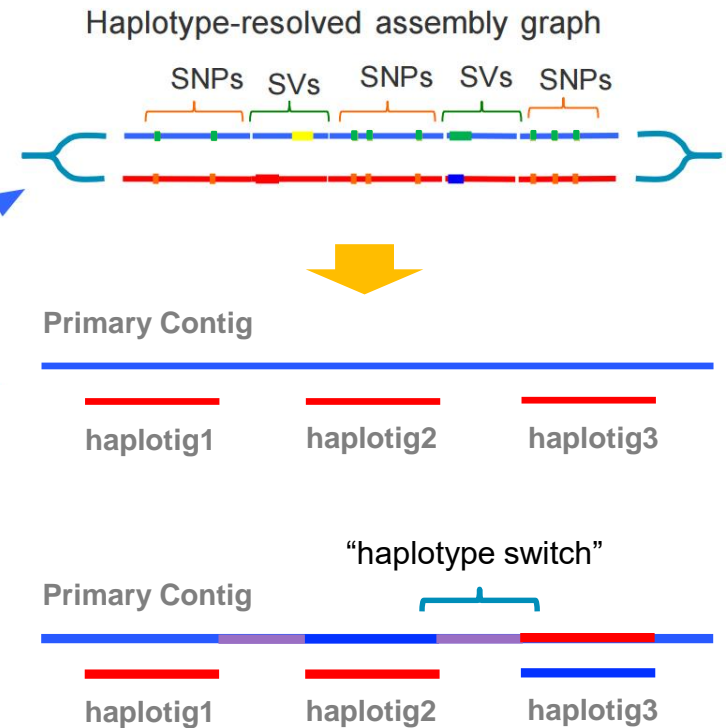
(a) FALCON



(b)



(c) FALCON-Unzip



HOW MUCH TO SEQUENCE FOR GENOME ANNOTATION?

EXAMPLE SEQUEL-SCALE EQUIVALENT YIELDS OF FL READS, GENES AND TRANSCRIPTS

| SPECIES | FL READS | GENES | TRANSCRIPTS | Using 3.0 Chemistry ⁺ : |
|----------------------------|-----------|--------|-------------|------------------------------------|
| Maize | 1,553,692 | 26,946 | 111,151 | ~4 SMRT Cell |
| Chicken (normalization) | 653,441 | 29,013 | 64,277 | ~1.5 SMRT Cell |
| Rabbit | 466,034 | 14,474 | 36,186 | ~1 SMRT Cell |
| R. necatrix | 330,373 | > 5000 | 10,616 | ~1 SMRT Cell |
| Zebra Finch | 405,736 | 7,228 | 17,437 | ~1 SMRT Cell |

⁺ based on 3.0 chemistry yield of ~400k FL reads per Sequel SMRT Cell

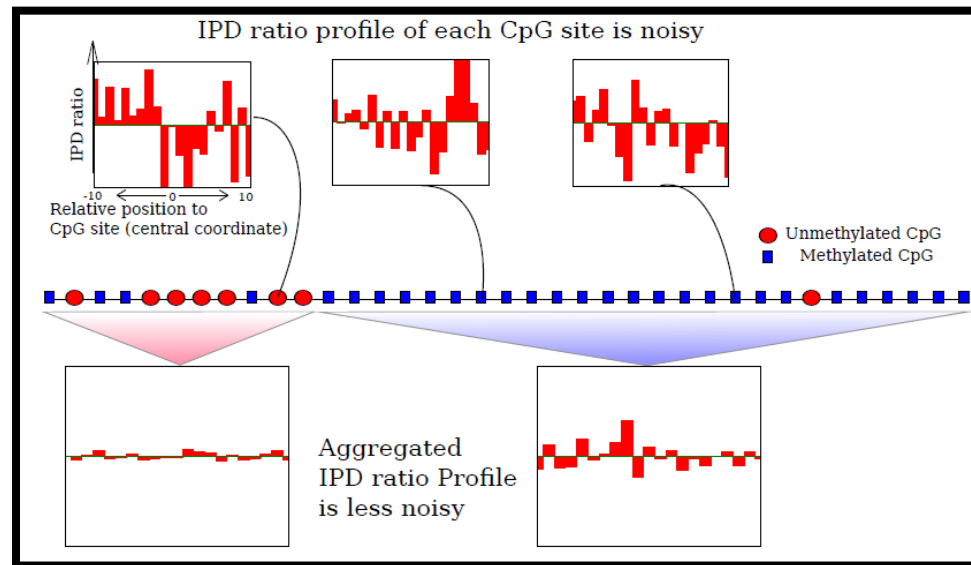
Wang et al., *Unveiling the complexity of the maize transcriptome by single-molecule long-read sequencing*, Nat Comm (2016)

Kuo et al., *Normalized long read RNA sequencing in chicken reveals transcriptome complexity similar to human*, BMC Genomics (2017)

Chen et al., *A transcriptome atlas of rabbit revealed by PacBio single-molecule long-read sequencing*, Sci Rep (2017)

Kim et al., *Characterization of the Rosellinia necatrix Transcriptome and Genes Related to Pathogenesis by Single-Molecule mRNA Sequencing*, Plant Patho J (2017)

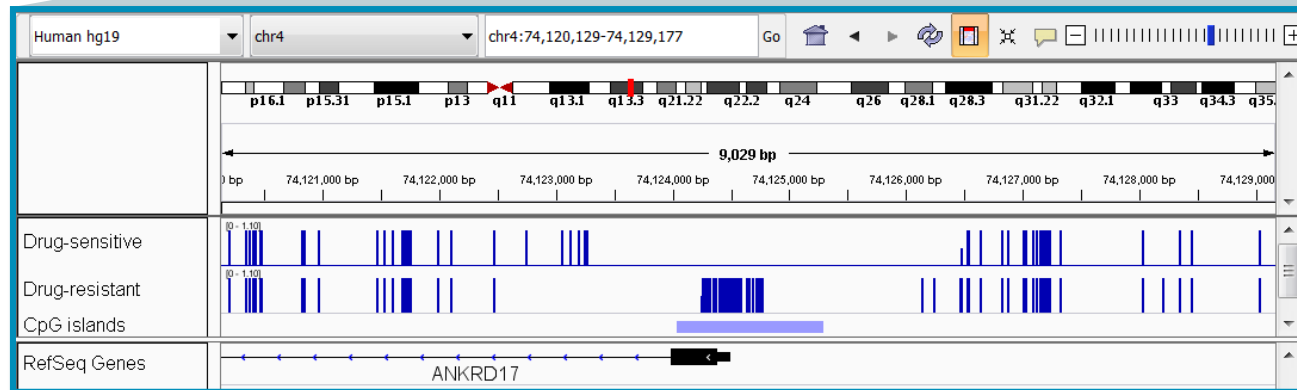
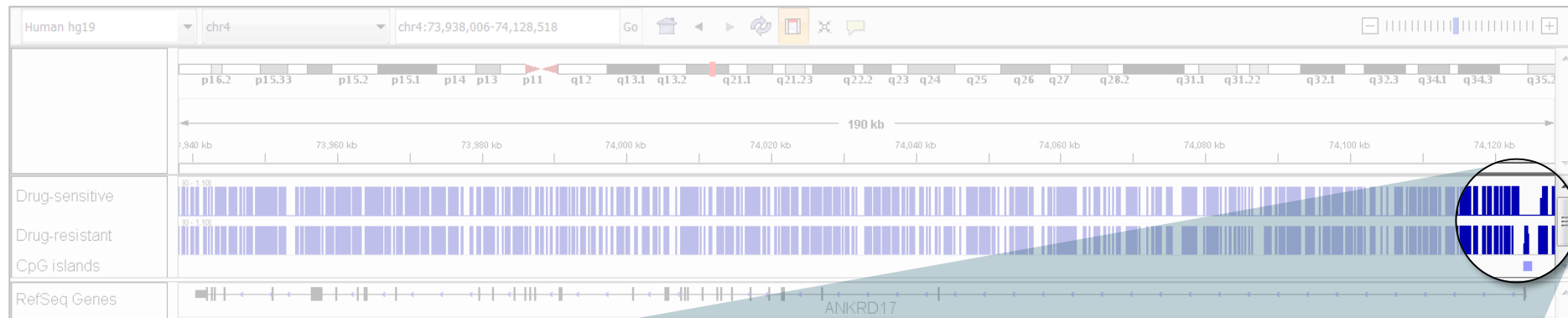
DETECTING HYPER- AND HYPO-METHYLATED CPG ISLANDS IN EUKARYOTIC DNA



- Signal strength for 5mC is weaker than for other methylated bases; requires high per-base coverage to detect
- Eukaryotic methylation is a regional phenomenon
- Prof. Shinichi Morishita (Univ. of Tokyo) developed a method to differentiate hypo- and hyper-methylated regions by integrating signals across CpG islands
- 16-fold per-strand coverage maximizes the accuracy of the results
- Algorithm is freely available at <https://github.com/hacone/AgIn>

EPIGENOME CHARACTERIZATION OF CANCER CELL LINES

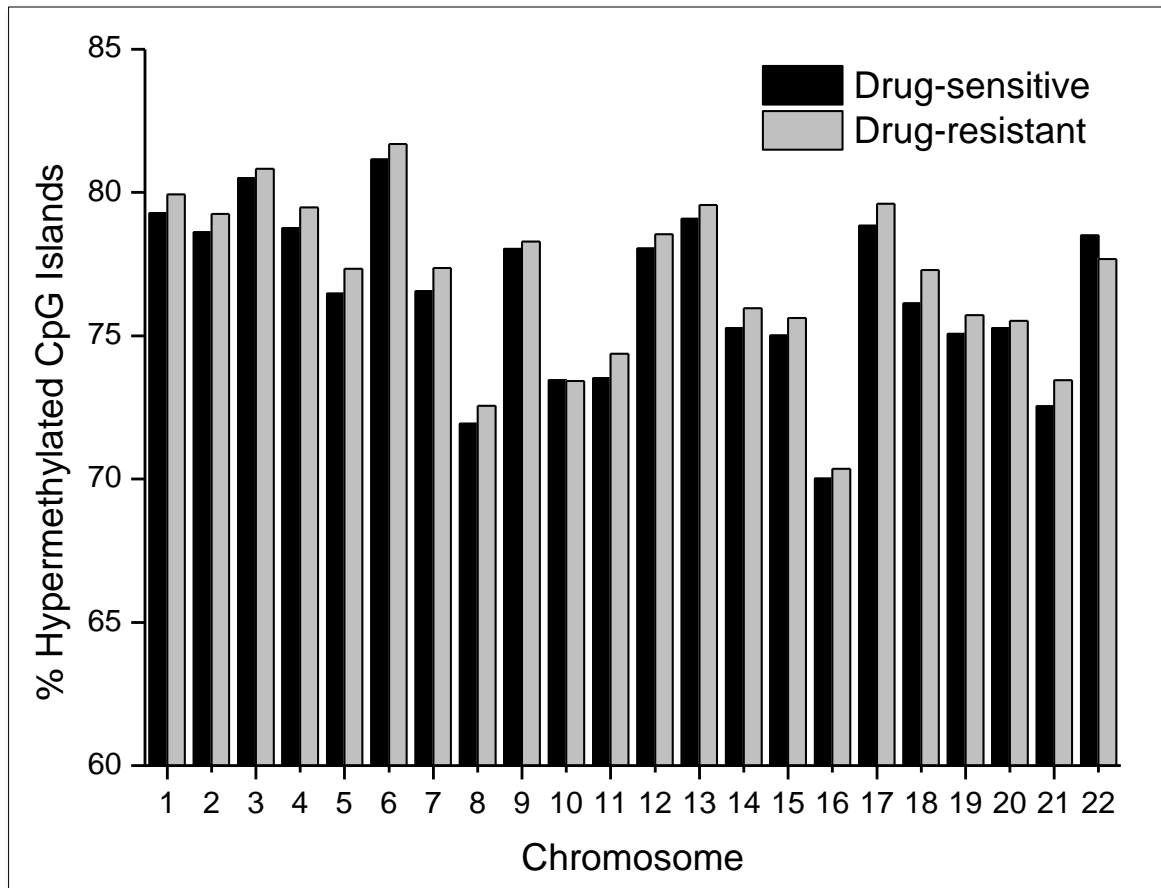
METHYLATION STATUS OF CpG ISLANDS ACROSS CHR4: *ANKRD17* (BREAST CANCER)



- Third-party custom analysis software available @ <https://github.com/hacone/AgIn>

EPIGENOME CHARACTERIZATION OF CANCER CELL LINE

Global methylation status:



More hypermethylated
CpG islands in drug-
resistant sample