

# **The Sequel II System**

An introduction to Single Molecule Real-Time (SMRT) Sequencing and its applications

Deborah Moine – Scientist II, Field Applications Support

For Research Use Only. Not for use in diagnostic procedures. © Copyright 2019 by Pacific Biosciences of California, Inc. All rights reserved.

## -SMRT Sequencing & Sequel II System

## - Applications :

- *de novo* assembly
- Variant detection
- Iso-Seq full-length cDNA sequencing
- Microbe characterization
- -Sequel II System updates

סאכן כל אכן כל איכ

#### **Long Reads**

- Tens of kilobases
- Sequence from 500 bp to >50,000 inserts

#### **High Accuracy**

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

#### **Single-Molecule Resolution**

- Sequence DNA or RNA
- Long reads with ≥Q20 (99%) single-molecule accuracy

#### **Uniform Coverage**

- No DNA amplification
- Least GC content and sequence complexity bias

#### **Simultaneous Epigenetic Detection**

- Characterize epigenome
- No separate sample preparation required

# SMRT Sequencing Advantages

#### אק כואכן כו איכן כו איכן כו איכ PACBIO\*

#### FROM SAMPLE TO SMRT SEQUENCING



millions of zero-mode waveguides (ZMWs)

sequence genomes, transcriptomes, and epigenomes

Prepare sequencing reaction

## ס-רכן כל ארכן כל ארכן

### SINGLE MOLECULE, REAL-TIME (SMRT) SEQUENCING



are measured in real time

SMRT Sequencing enables simultaneous collection of data from millions of wells using the natural process of DNA replication to sequencing long fragments of native DNA.



#### **HIGH CONSENSUS ACCURACY**

Achieves >99.999% (Q50)



Consensus accuracy is a function of coverage and chemistry. The data above is based on a bacterial genome run on the Sequel II System (1.0 Chemistry, Sequel II System Software v7.0). Single-molecule accuracy has similar coverage requirements.



#### **UNIFORM COVERAGE**



Mean coverage per GC window across a human sample. Data generated with a 35 kb human library on a Sequel II System using 1.0 Chemistry and Sequel II System Software v7.0

#### PACBIO PRODUCT RELEASES OVER THE LAST EIGHT YEARS

ליק כוי כן כוי כן כוי כן כוי כ

PACBIO\*



סאכן כל אכן כל איכ

### **SMRTBELL LIBRARY CONSTRUCTION**

# Existing Sequel library construction procedures are compatible with Sequel II System

- SMRTbell Template Prep Kit 1.0
- SMRTbell Express Template Prep Kit 2.0
  - Large insert gDNA library (≥10 kb)
  - Microbial multiplexing (10 kb)
  - Iso-Seq Express Method
  - Multiplexed amplicon sequencing
  - Low input DNA sequencing (150ng for a 300 Mb genome)
- Procedure & Checklist Preparing SMRTbell Libraries for HiFi Long Read Sequencing



Procedure & Checklist - Preparing SMRTbell<sup>®</sup> Libraries for HiFi Long Read Sequencing on Sequel<sup>®</sup> and Sequel II Systems

This document describes a method for constructing SMRTbell ibraries suitable for generating high accuracy long reads on the Sequel Systems. This procedure requires 16 ya of high molecular weight genomic DNA (gDNA). DNA is sheared to a mode of 15 kb using Diagenode's Megaruptor. SMRTbell libraries are prepared from sheared gDNA using PacBie's SMRTbell Template Prep Kit 1.0. SMRTbell libraries are then size-fractionated using Sage Science's SageELF. The fractions suitable for sequencing have insert sizes from 10 kb to 15 kb.

#### **Required Materials**

Item	Vendor
<u>gDNA QC</u> (one of the following) CHEF Mapper XA	Bio-Rad 170-3670 Sage Science PP10200
Pippin Pulse FEMTO Pulse	Agilent Technologies, Inc. (formerly Advanced Analytical Technologies, P-0003-0817)
DNA Quantitation Qubit 3.0 Fluorometer dsDNA HS Assay Kit	Life Technologies Q33216 Life Technologies Q32854
DNA Shearing Megaruptor Long Hydropores Hydrotubes	Diagenode B06010001 Diagenode E07010002 Diagenode C30010018
SMRTbell Library Preparation SMRTbell Template Prep kit 1.0 AMPure PB Beads Rotator 100% Ethanol, Molecular Biology Grade 2.0 mL DNA Lo-Bind Tubes Thermomixers	Pacific Biosciences 100-259-100 Pacific Biosciences 100-265-900 Any MLS Any MLS Eppendorf 022/31048 Any MLS
Fractionation SageELF System 0.75% Agarose Cassettes	Sage Science ELF0001 Sage Science ELD7510

Page 1

PN 101-714-400 Version 01 (April 2019)

**SEQUEL II SYSTEM – SEQUENCE WITH CONFIDENCE** 



#### Sequel System

סאכן כלא כין כלא כין כלא כין כלא כין כלא איר כין כלא ארמייט ארמייט איר איר פין כלא איר פין פי



1 million ZMWs SMRT Cell 1M

#### **Sequel II System**



8 million ZMWs SMRT Cell 8M

- ~8-fold increase in data yield
- Reduced project time
- Lower project cost
- Equivalent performance

**Increased Throughput Capacity With the Sequel II System** 

אק כוי כן כן יכן יכן יכ PACBIO\*

#### TWO MODES OF SMRT SEQUENCING

- Continuous Long Read Sequencing (CLR) - Circular Consensus Sequencing (CCS)



<u></u>
<u>AA\$\$74A\$4\$4\$4\$4\$\$</u>
<u>₩₩₽₽₩₩₩₽₩₽₩₽₩₽₩₽₩₽₩₩₽₩₩₩₽₩₩₽₩₩₽₩₽₩₽₩₽₩</u>
<u>,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,</u>
<u>₩₩₽₽₩₩₩₽₩₽₩₽₩₽₩₽₩₽₩₽₩₽₩₩₽₩₽₩₽₩₽₩₽₩₽₩₽₩</u>
<u></u>
₩₩₽₽₩₩₩₽₩₽₩₽₩₽₩₽₽₩₽₽₩₩₽₽₩₩₩₽₽₩₩₽₽₩₩₽₽₩

#### ¥¥\$₽¥¥¥\$¥\$¥\$¥\$¥**\$**¥**\$**\$₽₩¥\$₽¥¥¥\$¥\$¥\$¥\$¥\$¥\$<u>\$</u>\$¥¥\$**\$**\$¥¥\$**\$**\$¥¥\$¥\$¥\$¥\$¥\$¥\$¥\$¥

Consensus sequence



Subread 1	
,	<u></u>
,	
,	
,	
,	<u> </u>
-	<u></u>
,	<u></u>
Subread n	weentrevereneed een een ververe for the reverence of the

HiFi Read

Generate reads in ten's of kilobases

Generate high quality reads 1 – 20 kb

## סאק כלא כן כלא כן כלא כן כלא כן כלא כ

#### **TWO MODES OF SMRT SEQUENCING**

Continuous Long Read Sequencing (CLR)



CLR 1	<u>\************************************</u>
	*****************
,	<u></u>
	<u> </u>
	<u></u>
,	<u></u>
	<u></u>
,	<u></u>
	<u> </u>
	<u> </u>
	<u> </u>
CLR n	

Consensus sequence

- Circular Consensus Sequencing (CCS)



Subread 1	<u></u>
	······································
	<u>,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,</u>
	<u></u>
	¥8¥88¥8¥8¥88¥84849
	₩₩₽₽₩₩₩₽₩₽₩₽₩₽₩ <b>₽</b> ₽₽₩₩₽₽₩₩₩₽₩₽₩₽₩₽₩₽₩₽₩₽
	<u></u>
Subread n	₩₩₽₽₩₩₩₽₩₽₩₽₩₽₩ <mark>₽</mark> ₽₽₩₩₽₽₩₩₩₽₩₽₩₽₩₽₩₽₩₽₩₽

#### 

HiFi Read

Generate reads in ten's of kilobases

Generate high quality reads 1 – 20 kb

אר כן כל אכן כל ארכין כל ארכי



#### Read Length (bp)

Data shown above from a 35 kb size-selected *E. coli* library using the SMRTbell Template Prep Kit on a Sequel II System (1.0 Chemistry, Sequel II System Software v7.0, 15-hour movie). Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.

לא ליק כלי כן כל ייכן כל ייכן כל ייכ PACBIO Start with high-quality Circularized DNA double stranded DNA is sequenced in repeated passes in the second se \_\_\_\_\_\_ Ligate SMRTbell \_\_\_\_\_ adapters and size select The polymerase reads are trimmed of adapters to yield subreads 

Anneal primers and

bind DNA polymerase

Consensus is called from subreads

HIFI READ (99% accuracy with 4 passes)

## ארק כל אכן כל א

# HIFI READ PERFORMANCE



Data shown above from a 12 kb size-selected human library using the SMRTbell Template Prep Kit on a Sequel II System (1.0 Chemistry, Sequel II System Software v7.0, 30-hour movie). Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.



#### Read Length (bp)

Data shown above from a 11 kb size-selected human library using the SMRTbell Template Prep Kit on a Sequel II System (1.0 Chemistry, Sequel II System Software v7.0, 30-hour movie). Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.

ארק כל ארכן כל ארכין כל ארכי





Data shown above from a 11 kb size-selected human library using the SMRTbell Template Prep Kit on a Sequel II System (1.0 Chemistry, Sequel II System Software v7.0, 30-hour movie). Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.

THE SEQUEL II SYSTEM – SEQUENCE WITH CONFIDENCE

ס- רן כן כן ארכן כן כן ארכן כן כן ארכן כן כן ארכין כ



- ✓ Combines high accuracy and long read lengths
- Two sequencing modes for flexibility in your workflows



סער בן כל איכן כל איכ

#### **TWO MODES OF SMRT SEQUENCING**

### **Circular Consensus Sequencing (CCS) Mode**

Inserts 10-20 kb



## Continuous Long Read (CLR) Sequencing Mode

Inserts >25 kb, up to 175 kb



CLR 1	<u></u>
	<u></u>
	₩₩₽₽₩₩₩₽₩₽₩₽₩₽₩₽₩₽₩₽₩₩₽₩₩₽₩₩₽₩₩₽₩₩₽₩₩₽₩
	<u></u>
	<u>₩₽₩₽₽₩₽₩₽₩₽₽₽₽₩₽₽₩₽₩₽₩₽₩₽₩₽₩₽₩₽₽₩₩₽₽₩₩</u>
	≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈
CLR n	<u> </u>

#### 



Multi-molecule consensus sequence

WHOLE GENOME SEQUENCING (WGS) APPLICATIONS

סיק כו אכן כו אכן כו אכן כו אכן כו אכן





PACBIO\*

- Comprehensive detection of variants (SNVs, SVs, CNVs)
- High-quality, phased genome assembly

- Detection of structural variants (SVs)
- -Assembly of very large genomes



### WHAT CAN YOU DO WITH ONE SMRT CELL 8M?



PAC**BIO**®

HiFi Reads	Long Reads
Detect all variants – SNVs, indels, SVs, and CNVs – in a human genome (2 SMRT Cells 8M)	Detect structural and copy number variants in two human genomes
Assemble up to a 2 Gb genome	Assemble up to a 3 Gb genome

www.pacb.com/onesmrtcell

\*Read lengths, reads/data per SMRT Cell 8M, and other sequencing performance results vary based on sample quality/type and insert size. Prices, listed in USD, are approximate and may vary by region. Pricing includes library and sequencing reagents run on a Sequel II System and does not include instrument amortization or other reagents.



# **De novo Assembly**

Using Whole Genome Sequencing on the Sequel II System

סאק כלא כן כל איכן כל איכן כל איכן כל איכ מין כל איכ פ

#### **ANSWERING: WHAT'S THE SAME AND WHAT'S DIFFERENT?**

High-quality genome assemblies are the most comprehensive way to assess genomic diversity among and between species.



Crop Improvements



Conservation



Population or Disease-specific Studies



Animal Health & Breeding



**Disease Prevention** 

ס- כן ארמין כן אימין פרא פרי כ

#### **DRAFT VS HIGH-QUALITY GENOME ASSEMBLIES**

You see more, and therefore have more opportunities to see differences, with a high-quality genome assembly.



סאכן כלא כין כלא כין כלא כין כלא כין כלא כי

### DRAFT VS HIGH-QUALITY GENOME ASSEMBLIES

You see more, and therefore have more opportunities to see differences, with a high-quality genome assembly.

#### Contig:



סאכן כל אכן כל איכ

### UNIQUE CHALLENGES OF GENOME ASSEMBLY

#### **Size and Complexity**

- Human genome has over 3 billion base pairs
- Plants often have even larger genomes

#### **Extreme Repeat Content**

- Maize >60%
- Wheat >80%

#### **Each Project is Unique**

- Ranges in size, ploidy, and repeat content
- Custom strategy is commonly needed



ס-רק כל ארכן כ

#### WHY DO LONG READS MATTER?



Long, accurate reads span repetitive elements and allow assembly of repeats that are very similar Even really **f** big genomes



# CHOOSE THE SEQUENCING MODE THAT'S RIGHT FOR YOUR PROJECT



- Highly accurate long reads with minimum accuracy of Q20 (99%)
- Small file sizes and fast analysis time
- Assemble up to a 2 Gb genome in a single SMRT Cell 8M\*
- Run up to 200 samples (2 Gb) per year, per system\*



PACBIO

- Longest reads, with half of data
  >50 kb and maximum read lengths up to 175 kb
- Short sequencing run times
- Assemble up to a 3 Gb genome in a single SMRT Cell 8M\*
- Run up to 400 samples (3 Gb) per year, per system\*

\*Read lengths, reads/data per SMRT Cell 8M, and other sequencing performance results vary based on sample quality/type and insert size. Prices, listed in USD, are approximate and may vary by region. Pricing includes library and sequencing reagents run on a Sequel II System and does not include instrument amortization or other reagents.

WORKFLOWS: FROM DNA TO COMPREHENSIVE GENOME ASSEMBLIES

# HIFI HIFI Reads

>15 µg of unamplified genomic DNA

Express Template Prep Kit 2.0 Size select 15-20 kb inserts

Use Circular Consensus Sequencing (CCS) mode

1			
	_	_	н
	_	_	н
	_	_	н
L			

**LIBRARY PREP** 

SMRT SEOUENCING Use Continuous Long Read (CLR) sequencing mode

 $>5 \mu g$  of unamplified genomic DNA

Express Template Prep Kit 2.0

Size select >30 kb inserts

PAC**BIO**®

Generate HiFi reads with CCS analysis

Assemble and phase HiFi reads with FALCON and FALCON-Unzip3

DATA ANALYSIS

Assemble and phase long reads with PacBio analytical portfolio:

- -SMRT Analysis
- -FALCON and FALCON-Unzip
- -Bioinformatics service providers



## סאק כלא כן כלא כן כלא כן כלא כן כלא כ

### **GENERATE CONTIGUOUS GENOME ASSEMBLIES**

Dataset	Human		Rice	
Mode	Long Reads	HiFi Reads	Long Reads	HiFi Reads
Size Selection	BP >15 kb	15 kb ELF	BP >30 kb	17 kb ELF
Coverage	50-fold	22-fold	60-fold	20-fold
# of SMRT Cells 8M	2	3	1	1
Contig N50 (Mb)	12.6	30.5	11.2	10.7

Megabase size contig N50s

#### GENERATE COMPLETE AND ACCURATE GENOME ASSEMBLIES

🐼 ליק כוי כן כויי כן כויי כ

PACBIO°



Accuracies >Q40 (99.99%) >94% of genes in frame סאכן כלא כין כלא כין כלא כין כלא אר כין כ

### **REDUCED ANALYSIS TIME WITH HIFI READS**



#### Compute times for *de novo* assembly of a human genome

Data Type		HiFi Reads	Long Reads
Input File Type		CCS.FASTQ.GZ	SUBREADS.BAM
Input File Size (GB)		48	323
Read Correction Method		CCS Analysis	Pre-assembly
Time to Results C (Hours) C	Read Correction	17.5	43.5
	Contig Assembly	13.7	18.9
Analyses run with PacBio recor	mmended compute infrastruc	~31 hrs	~62 hrs

## Fast analysis time



# Variant Detection

Using Whole Genome Sequencing on the Sequel II System

סאק כל אכן כל אכן כל אכן כל אכן כל ארכן כל ארכ

### **TYPES OF VARIANTS IN A GENOME**

SMRT Sequencing provides comprehensive detection of all variant types.



#### A COMPREHENSIVE VIEW OF THE GENOME

SMRT Sequencing provides a view into *all* variation between two humans.



### A COMPREHENSIVE VIEW OF THE GENOME

SMRT Sequencing provides a view into *all* variation between two humans.


סאכן כלאכן כלא ארכ

#### A COMPREHENSIVE VIEW OF THE GENOME

SMRT Sequencing provides even coverage across difficult to sequence regions of the genome.



Almost no coverage with prior tech

PacBio reads sequence straight through and detect variants, some falling in coding regions

#### A COMPREHENSIVE VIEW OF THE GENOME

SMRT Sequencing provides a view into *all* variation between two humans.



Long insertions Events in repeat regions 

#### A COMPREHENSIVE VIEW OF THE GENOME

SMRT Sequencing provides long read lengths to span large structural variants.



### **COMPREHENSIVE VARIANT DETECTION WITH HIFI READS**

סיק כרי כן כרי כן כרי כן כרי כ

#### **The HiFi Difference**

- High precision and recall for SNVs, indels, and SVs
- Detect 5% more variants in "medical exome"
- Detect structural variants
- Phase variants into haplotypes

Aaron M. Wenger <sup>© 134</sup> , Paul Peluso <sup>134</sup> , William J. Rowell <sup>©1</sup> , Pi-Chuan Chang <sup>©2</sup> , Richard J. Hall <sup>1</sup> , Gregory T. Concepcion <sup>©1</sup> , Jana Ebler <sup>24,5</sup> , Arkarachai Fungtammasan <sup>6</sup> , Alexey Kolesnikov <sup>2</sup> , Nathan D. Olson <sup>©7</sup> , Armin Töpfer <sup>1</sup> , Michael Alonge <sup>6</sup> , Medhat Mahmoud <sup>9</sup> , Yufeng Qian <sup>1</sup> , Chen-Shan Chin <sup>©6</sup> , Adam M. Phillipp <sup>10</sup> , Michael C. Schatz <sup>8</sup> , Gene Myers <sup>11</sup> , Mark A. DePristo <sup>2</sup> , Jue Ruan <sup>©1</sup> , Tobias Marschall <sup>© 34</sup> , Fritz J. Sedlazeck <sup>©1</sup> , Justin M. Zook <sup>©7</sup> , Heng Li <sup>01</sup> , Sergey Koren <sup>10</sup> , Andrew Carroll <sup>2</sup> , David R. Rank <sup>© 1*</sup> and Michael W. Hunkapiller <sup>©1*</sup> The DNA sequencing technologies in use today produce either highly accurate short reads or less-accurate long reads. report the optimization of circular consensus sequencing (CCS) to Improve the accuracy of single-moleculer lead-time (SM) sequencing (Patibio) and generate highly accurate (99.8%) (Ding high-fidelity (Hifi) reads with an average length (13.5) bases (kb). We applied our approach to sequence the well-characterized human HGO02/NA24385 genome and obtained- cion and recall rates of at least 99.9% for single-nucleotide variants (SN+9), 59.5% for insert <sup>1</sup> (Indels) and 95.99% for structural variants. We estimate that 2.434 discordances or exceeds the ability and Variants and structural variants. We action ac produce a core	improves var	iant detection and assembly of a
Aaron M. Wenger <sup>© 1,4</sup> , Paul Peluso <sup>1,4</sup> , William J. Rowell <sup>©</sup> <sup>1</sup> , Pi-Chuan Chang <sup>©<sup>2</sup></sup> , Richard J. Hall <sup>1</sup> , Gregory T. Concepcion <sup>©</sup> <sup>1</sup> , Jana Ebler <sup>2,4,5</sup> , Arkarachai Fungtammasan <sup>6</sup> , Alexey Kolesnikov <sup>2</sup> , Nathan D. Olson <sup>©<sup>2</sup></sup> , Armin Töpfer <sup>1</sup> , Michael Alonge <sup>8</sup> , Medhat Mahmoud <sup>9</sup> , Yufeng Qian <sup>1</sup> , Chen-Shan Chin <sup>©<sup>4</sup></sup> , Adam M. Phillippy <sup>10</sup> , Michael C. Schatz <sup>2</sup> , Gene Myers <sup>10</sup> , Mark A. DePristo <sup>2</sup> , Jue Ruan <sup>©<sup>2</sup></sup> , Tobias Marschall <sup>© 3,4</sup> , Fritz J. Sedlazeck <sup>©<sup>9</sup></sup> , Justin M. Zook <sup>©<sup>7</sup></sup> , Heng Li <sup>© 1,8</sup> Sergey Koren <sup>10</sup> , Andrew Carroll <sup>2</sup> , David R. Rank <sup>© 1*</sup> and Michael W. Hunkapiller <sup>© 1*</sup> The DNA sequencing technologies in use today produce either highly accurate short reads or less-accurate long reads. report the optimization of circular consensus sequencing (CSS) to improve the accuracy of single-molecular telepishot of 3.51 bases (kb). We applied our approach to sequence the well-characterized human HG002/NA24385 genome and obtained- ciaion and recall rates of at least 99.91% for single-nucleotide variants (SNV), 95.98% for insert (indels) and 95.99% for structural variants. We estimate that 2.424 discordances at the bil <sup>11</sup> singla variants and structural variants. We acho matches or exceeds the ability isonal variants and structural variants. We also a discource of the application of the sequence of the structural variants the structural variants the structural variants. We structural variants constructural variants constructural variants and base or exceeds the ability indiciab and 95.09% for structural variants. We also and more and estimate the structural variants	human genor	ne
Gregory T. Concepcion <sup>®</sup> , Jana Ebler <sup>3,45</sup> , Arkarachai Fungtammasan <sup>6</sup> , Alexey Kolesnikov <sup>2</sup> , Nathan D. Olson <sup>®</sup> , Armin Töpfer <sup>1</sup> , Michael Alonge <sup>8</sup> , Medhat Mahmoud <sup>9</sup> , Yufeng Qian <sup>1</sup> , Chen-Shan Chin <sup>®</sup> , Adam M. Phillipp <sup>10</sup> , Michael C. Schatz <sup>8</sup> , Gene Myers <sup>11</sup> , Mark A. DePristo <sup>2</sup> , Jue Ruan <sup>® 12</sup> , Tobias Marschall <sup>® 34</sup> , Fritz J. Sedlazeck <sup>®</sup> , Justin M. Zook <sup>®</sup> , Heng Li <sup>®</sup> , Sergey Koren <sup>10</sup> , Andrew Carroll <sup>2</sup> , David R. Rank <sup>® 1*</sup> and Michael W. Hunkapiller <sup>® 1*</sup> The DNA sequencing technologies in use today produce either highly accurate short reads or less-accurate long reads. sequencing (Patibio) and generate highly accurate (99.8%) (Ing) high-fidelity (Hif) reads with an average length of 13.51 bases (kb). We applied our approach to sequence the well-characterized human HGO02/NA24385 genome and obtained- cion and recuir lrates of at least 99.9% for single-nucleotide variants (SNV), 95.98% for inset <sup>11</sup> (Indels) and 95.99% for structural variants. We estimate that 24.34 discordances are berrow zenome assembly using CCS reads alone produced a core	Aaron M. Wenger <sup>(2),14</sup> , Pa	ul Peluso <sup>1,14</sup> , William J. Rowell <sup>(3)</sup> , Pi-Chuan Chang <sup>(3)</sup> , Richard J. Hall <sup>1</sup> ,
Nathan D. Olson <sup>©</sup> , Armin Töpter, Michael Alonge <sup>°</sup> , Medhat Mahmoud <sup>°</sup> , Yuteng Qian, Chen-Shan Chin <sup>©</sup> <sup>6</sup> , Adam M. Phillippy <sup>®</sup> , Michael C. Schatz <sup>8</sup> , Gene Myers <sup>11</sup> , Mark A. DePristo <sup>2</sup> , Jue Ruan <sup>©</sup> <sup>1</sup> , Tobias Marschall <sup>©,14</sup> , Fritz J. Sedlazeck <sup>©,1</sup> , Justin M. Zook <sup>©</sup> , Heng Li <sup>©,14</sup> , Sergey Koren <sup>10</sup> , Andrew Carroll <sup>2</sup> , David R. Rank <sup>©,1*</sup> and Michael W. Hunkapiller <sup>©,1*</sup> The DNA sequencing technologies in use today produce either highly accurate short reads or less-accurate long reads. sequencing (Patibio) and generate highly accurate (99.8%) (Inon jhigh-fidelity (Hif) reads with an average length of 13.51 bases (kb). We applied our approach to sequence the well-characterized human HGO02/NA24385 genome and obtained- cion and recall rates of at least 99.9% for single-nucleotide variants (SNV), 95.98% for inset <sup>14</sup> (Indels) and 95.99% for structural variants. We sum that 24.44 discordances are berrow zenome assembly using CCS reads alone produced a cort	Gregory T. Concepcion <sup>(0)</sup> ,	Jana Ebler <sup>3,4,5</sup> , Arkarachai Fungtammasan <sup>6</sup> , Alexey Kolesnikov <sup>2</sup> ,
Chen-Shan Chin <sup>®</sup> , Adam M. Phillippy <sup>®</sup> , Michael C. Schatz <sup>*</sup> , Gene Myers <sup>*</sup> , Mark A. DePristo <sup>*</sup> , Jue Ruan <sup>®</sup> <sup>1</sup> , Tobias Marschall <sup>®34</sup> , Fritz J. Sedlazeck <sup>®9</sup> , Justin M. Zook <sup>®7</sup> , Heng Li <sup>®1</sup> , Sergey Koren <sup>®</sup> , Andrew Carroll <sup>2</sup> , David R. Rank <sup>®1*</sup> and Michael W. Hunkapiller <sup>®1*</sup> The DNA sequencing technologies in use today produce either highly accurate short reads or less-accurate long reads. report the optimization of circular consensus sequencing (CCS) to improve the accuracy of single-molecular lead-time (SM sequencing (Pacific) and generate highly accurate (99.8%) (Ino) high-fieldity (Hif) reads with an average length of 13.51 bases (kb). We applied our approach to sequence the well-characterized human HG002/NA24385 genome and obtained- ciaion and recil rates of at least 99.91% for single-nucleative variants (SNV), 95.98% for insert <sup>1</sup> (indels) and 95.99% for structural variants. We structure that discontances or exceeds the ability of the sequence of the structural variants that 24.44 discontances are exceeded by a structural variants. We structure the accurate and structural variants. We school and produce accurate constructural variants can be phased to Benovo senome assembly using CCS reads alone produced a constructural variants.	Nathan D. Olson 0', Armir	1 Topfer', Michael Alonge°, Medhat Mahmoud°, Yufeng Qian',
Jue Ruan 0 <sup>14</sup> , Iobias Marschall 0 <sup>14</sup> , Fritz J. Sediazeck 0 <sup>4</sup> , Justin M. Zook 0 <sup>4</sup> , Heng Li 0 <sup>14</sup> , Sergey Koren <sup>10</sup> , Andrew Carroll <sup>2</sup> , David R. Rank 0 <sup>14</sup> and Michael W. Hunkapiller 0 <sup>14</sup> The DNA sequencing technologies in use today produce either highly accurate short reads or less-accurate long reads. report the optimization of circular consensus sequencing (CCS) to improve the accuracy of single-molecule read-time (SM sequencing (PacBio) and generate highly accurate (99.8%) long high-fidelity (HiF) reads with an average length of 13.51 bases (kb). We applied our approach to sequence the well-characterized human HG002/NA24385 genome and obtained cision and recall rates of at least 99.9% for single-nucleatile variants (SNN), 95.98% for insert (Indels) and 95.99% for structural variants. We estimate that 2,434 discordances are benovo senome assembly using CCS reads alone produced a corr	Chen-Shan Chin@®, Adam	M. Phillippy <sup>10</sup> , Michael C. Schatz <sup>8</sup> , Gene Myers <sup>11</sup> , Mark A. DePristo <sup>2</sup> ,
Sergey Koren <sup>™</sup> , Andrew Carroll <sup>*</sup> , David R. Rank <sup>©</sup> <sup>™</sup> and Michael W. Hunkapiller <sup>©</sup> <sup>™</sup> The DNA sequencing technologies in use today produce either highly accurate short reads or test-accurate long reads. report the optimization of circular consensus sequencing (CCS) to improve the accuracy of single-molecule real-time (SM sequencing (PaciBio) and generate highly accurate (99.8%) long high-fidelity (HiT) reads with an average length of 13.51 bases (tob). We applied our approach to sequence the well-characterized human HG002/NA24385 genome and obtained- cision and recall rates of at least 99.9% for single-nucleatile variants (SMVs), 95.98% for insert (indels) and 95.99% for structural variants. We estimate that 24.34 discordances are beild (GAB) benchmark set. Nearly all (99.64%) variants can be phased the Gnew second essembly using CCS reads alone produced a con-	Jue Ruan 0 ", Tobias Mars	chall <sup>103</sup> , Fritz J. Sedlazeck <sup>107</sup> , Justin M. Zook <sup>107</sup> , Heng Li <sup>101</sup> ,
The DNA sequencing technologies in use today produce either highly accurate short reads or less-accurate long reads. report the optimization of circular consensus sequencing (CCS) to improve the accuracy of single-molecular earlieme (SM sequencing (PaciBio) and generate highly accurate (99.8%) long high-fieldity (Hiff) reads with an average length of 13.51 bases (kb). We applied our approach to sequence the well-characterized human HG002/NA24385 genome and obtained cision and recall rates of at least 99.9% for single-nucleotide variants (SNVs), 95.98% for insert- (indels) and 95.99% for structural variants. Our CCS method matches or exceeds the ability of insert- iand variants and structural variants. We estimate that 2.434 discordances are ble' (GIAB) benchmark set. Nearly all (99.64%) variants can be phased.	Sergey Koren <sup>10</sup> , Andrew Ca	arroll', David R. Rank 😳 ** and Michael W. Hunkapiller 🙂 **
cision and recall rates of at least 99.91% for single-nucleotide variants (SNVs), 95.98% for insert (indels) and 95.99% for structural variants. our CCS method matches or receeds the ability small variants and structural variants. We estimate that 2,434 discordances are the (GIAB) benchmark set. Nearly all (99.64%) variants can be phased Denvo sceneme assembly using CCS reads alone produced a continue of the structural structural variants and structural variants and structural variants.	The DNA sequencing technolog report the optimization of circul sequencing (PacBio) and genera bases (kb). We applied our appro-	ies in use today produce either highly accurate short reads or less-accurate long reads. V ar consensus sequencing (CCS) to improve the accuracy of single-molecule real-time (SMR te highly accurate (99.8%) long high-fidelity (Hifi) reads with an average length of 13.5 kil aoch to sequence the well-characterized human HGO02/NA24385 genome and obtained
small variants and structural variants. We estimate that 2,434 discordances and the structural variants. We estimate that 2,434 discordances are structured by the structure of	cision and recall rates of at least (indels) and 95,99% for structur	t 99.91% for single-nucleotide variants (SNVs), 95.98% for insertional sector insertion of the sector of the secto
tle' (GIAB) benchmark set. Nearly all (99.64%) variants can be phased Denovo genome assembly using CCS reads alone produced a communication	small variants and structural var	iants. We estimate that 2,434 discordances are
	tie' (GIAB) benchmark set. Near	ly all (99.64%) variants can be phase

nature

PAC**BIO**®

A D T I C I E S

Wenger, A. et al., <u>Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome</u>. 2019. *Nature Biotechnology.* 

## MEASURING PRECISION AND RECALL – GENOME IN A BOTTLE BENCHMARK

Reference



## Well-characterized benchmark variants

PACBIO\*



#### **MORE VARIANTS IN MEDICALLY-RELEVANT GENES**

איק כויכן כויכן כן כויכן כי כ

PACBIO\*

% problem exons resolved	• ORIGINAL RESEARCH ARTICLE	Genes	8						
100%	ABCC6, ABCD1, ACAN, ACSM2B, AKR1C2, ALG1, ANKRD11, BCR, CATSPER2, CD177, CEL, CES1, CFH, CFHR1, CFHR3, CFHR4, CGB, CHEK2, CISD2, CLCNKA, CLCNKB, CORO1A, COX10, CRYBB2, CSH1, CYP11B1, CYP11B2, CYP21A2, CYP2A6, CYP2D6, CYP2F1, CYP4A22, DDX11, DHRS4L1, DIS3L2, DND1, DPY19L2, DUOX2, ESRRA, F8, FAM120A, FAM205A, FANCD2, FCGR1A, FCGR2A, FCGR3A, FCGR3B, FLG, FLNC, FOXD4, FOXO3, FUT3, GBA, GFRA2, GON4L, GRM5, GSTM1, GYPA, GYPB, GYPE, HBA1, HBA2, HBG1, HBG2, HP, HS6ST1, IDS, IFT122, IKBKG, IL9R, KIR2DL1, KIR2DL3, KMT2C, KRT17, KRT6A, KRT6B, KRT6C, KRT81, KRT86, LEFTY2, LPA, MST1, MUC5B, MYH6, MYH7, NEB, NLGN4X, NLGN4Y, NOS2, NOTCH2, NXF5, OPN1LW, OR2T5, OR51A2, PCDH11X, PCDHB4, PGAM1, PHC1, PIK3CA, PKD1, PLA2G10, PLEKHM1, PLG, PMS2, PRB1, PRDM9, PROS1, RAB40AL, RALGAPA1, RANBP2, RHCE, RHD, RHPN2, ROCK1, SAA1, SDHA, SDHC, SFTPA1, SFTPA2, SIGLEC14, SLC6A8, SMG1, SPATA31C1, SPTLC1, SRGAP2, SSX7, STAT5B, STK19, STRC, SULT1A1, SUZ12, TBX20, TCEB3C, TLR1, TLR6, TMEM231, TNXB, TRIOBP, TRPA1, TTN, TUBA1A, TUBB2B, UGT1A5, UGT2B15, UGT2B17, UNC93B1, VCY, VWF, WDR72, ZNF419, ZNF592, ZNF674								
[75%, 100%)	ANAPC1, C4A, C4B, CHRNA7, CR1, DUX4, FCGR2B, HYDIN, OTOA, PDPK1, TMLHE								
[50%, 75%)	ADAMTSL2, CDY2A, DAZ1, GTF2I, NAIP, OCLN, RPS17								
[25%, 50%)	DAZ2, DAZ3, KIR3DL1, OPN1MW, PPIP5K1		_ 7						
(0%, 25%)	NCF1, RBMY1A1	11	/ '						
0%	BPY2, CCL3L1, CCL4L1, CDY1, CFC1, CFC1B, GTF2IRD2, HSFY1, MRC1, OR4F5, PRY, PRY2, SMN1, SMN2, TSPY1, XKRY	16							

Mandelker, D. et al., Navigating highly homologous genes in a molecular diagnostic setting: a resource for clinical next-generation sequencing. 2016. Genetics in Medicine 18, 1282-1289 Wenger, A. et al., Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. 2019. Nature Biotechnology.

#### **15-FOLD HIFI READ COVERAGE RECOMMENDATION**

ס- כן כן כן כן כן כן כן כן כן ארכן כן כ



Wenger, A. et al., Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. 2019. Nature Biotechnology.

סיק כוי כן כוי כן כוי כן כוי כ

#### WORKFLOW: FROM DNA TO COMPREHENSIVE VARIANT DETECTION WITH HIFI READS



PACBIO\*



3-6 µg DNA input per sample

Prepare 1 SMRTbell library per sample with Express Template Prep Kit 2.0 (15-20 kb size fraction)



SMRT SEOUENCING Sequence 2 SMRT Cells 8M per sample for 30 hrs (in CCS mode)



Call small variants with DeepVariant or GATK Call structural variants with pbsv

DATA ANALYSIS

סאק כלא כן כל איכן כל איכן כל איכן כל איכ פין כל איכ פין כל איכ

#### **STRUCTURAL VARIANT DETECTION IN A SINGLE SMRT CELL 8M**

SMRT Sequencing detects novel SVs of all types and lengths with base pair resolution of breakpoints.



long insertions events in repeat regions סאק כל אכן כל אכן כל אכן כל אכן כל אכן כל אכן כל ארכן כל איכן כל איכ

#### STRUCTURAL VARIANT DETECTION WITH LONG READS

#### **The SMRT Sequencing Difference**

High precision and recall for detection of:

- Insertions and deletions
- Inversions
- Translocations
- Copy-number variants

Experiment type	SMRT Sequencing	+ Short-read Sequencing
Disease Research	10-fold CLR coverage to detect SVs (multiplex 2 samples per SMRT Cell 8M)	30-fold coverage to detect SNVs and indels
Population Genetics	5- to 10-fold CLR coverage per sample to discover SVs (multiplex 2-3 samples per SMRT Cell 8M)	30-fold coverage per sample in a large cohort to genotype discovered SVs

#### סאכן כל אכן כל איכן כל איכן כל איכ

#### **CASE STUDIES**

Journal of Human Genetics	Cell
ARTICLE	Characterizing the Major Structural Variant Alleles of the Human Genome
A 12-kb structural variation in progressive myoclonic epilepsy was newly identified by long-read whole-genome sequencing Takeshi Mizuguchi <sup>1</sup> · Takeshi Suzuki <sup>2</sup> · Chihiro Abe <sup>2</sup> · Ayako Umemura <sup>2</sup> · Katsushi Tokunaga <sup>3</sup> · Yosuke Kawai <sup>3</sup> · Minoru Nakamura <sup>4</sup> · Masao Nagasaki <sup>5</sup> · Kengo Kinoshita <sup>6,7,8</sup> · Yasunobu Okamura <sup>6,7</sup> · Satoko Miyatake <sup>1,9</sup> · Noriko Miyake <sup>1</sup> · Naomichi Matsumoto <sup>1</sup> Received: 7 December 2018 / Revised: 22 January 2019 / Accepted: 22 January 2019 o The Author(s) under exclusive licence to The Japan Sodety of Human Genetics 2019	Graphical Abstract Authors Peter A. Audano, Arvis Sulovari, Tina A. Graves-Lindsay,, Yang I. Li, Richard K. Wilson, Evan E. Eichler Correct reference assembly Patch GRCh38 Align short reads
Abstract	GRCh38



human genome.

Mizuguchi T, et al. <u>A 12-kb structural variation in progressive myoclonic epilepsy was newly identified by long-read whole-genome sequencing</u>. 2019. *J* Human Genetics, 64, 359-368.

Audano P, et al. Characterizing the Major Structural Variant Alleles of the Human Genome. 2019. Cell, 176 (3), 663-675.

Roach, M.J., et al. Population sequencing reveals clonal diversity and ancestral inbreeding in the grapevine cultivar Chardonnay. 2018. PLOS Genetics 14.

### סיק כוא כן כוא כן כוא כן כוא כ

PACBIO\*

#### **RECALL VS COVERAGE FOR SV DETECTION**



Precision and recall for a single SMRT Cell 8M is for a human-sized genome. For a genome >3 Gb, more SMRT Cells may be required. WORKFLOW: FROM DNA TO SENSITIVE STRUCTURAL VARIANT DETECTION

ליכן כו יכן כו יכן כו יכן כו יכן



PACBIO



Prepare 3-6 µg DNA (>20 kb) input per sample

Prepare 1 SMRTbell library per sample using barcoded adapters and Express Template Prep Kit 2.0 (size select >15 kb)



SMRT SEOUENCING Multiplex 2 samples per 1 SMRT Cell 8M, sequencing for 15 hrs (in CLR mode)



Call structural variants jointly with pbsv

DATA ANALYSIS



## RNA Sequencing with the Iso-Seq Method for Transcriptome Analysis



#### **DETERMINATION OF TRANSCRIPT ISOFORMS**



Full-length cDNA Sequence Reads Splice Isoform Certainty – <u>No Assembly Required</u> סאק כל אכן כל איכן כל איכן כל איכן כל איכן כל א

#### **ISO-SEQ METHOD: FULL-LENGTH RNA SEQUENCING**

Sequence full-length cDNA sequences – from 5' end to the poly-A tail – without the need for transcript reconstruction.

The PacBio Iso-Seq method allows you to:

- Profile whole transcriptomes exhaustively at the isoform level
- Discover novel genes and isoforms
- Identify exon skipping events and alternative 5' / 3' sites
- Characterize function without reference genomes
- Combine with and complement RNA-seq to quantify at isoform-level



#### **IMPROVED GENOME ANNOTATION WITH ISO-SEQ**

לא ליק כלי כן כלי כן כלי כן כלי כ



PACBIO\*

Plutella xylostella Diamondback moth



Alternative splicing contributes to a **much more diverse transcript set** compared to original gene model

Zhao, Q. et al., 2019. Genome-wide profiling of the alternative splicing provides insights into development in Plutella xylostella, BMC Genomics, 20:463.

#### **SEQUENCING THE CANCER GENOME & TRANSCRIPTOME**

ליק כוי כן יכן יכן יכן יכ





PAC**BIO**°

Supplementary Figure 18. Ribbon plot of "3-hop" KLHDC2-SNTB1 gene fusion captured by long reads. This is a "3-hop" gene fusion in SK-BR-3 created by a series of three variants (A). These variants are captured together in several individual SMRT sequencing reads, one of which is shown in (B).

## In total, Iso-Seq identified **15 fusion genes** with genomic evidence



Further reading: Medium blogpost

Nattestad, M. et al. Complex rearrangements and oncogene amplifications revealed by long-read DNA and RNA sequencing of a breast cancer cell line. Genome Res. 1–19 (2018).

#### WORKFLOW: THE ISO-SEQ METHOD FOR TRANSCRIPTOME CHARACTERIZATION



300 ng of total RNA input per sample Prepare 1 SMRTbell library per sample with Express Template Prep Kit 2.0



Use the Sequel II System to generate up to 4 million fulllength, non-concatemer (FLNC) reads per SMRT Cell 8M

SMRT SEQUENCING



DATA ANALYSIS

Use the Iso-Seq analysis in SMRT Link to output high-quality, full-length transcript sequences

PRONEX BEAD PURIFICATION ENABLES MODULATION OF FULL-LENGTH TRANSCRIPT SIZE-DISTRIBUTION



סיכן כל ארכן כל

סאק כלא כין כל איכין כל איכין כל איכין כל איכין אי

#### ISO-SEQ EXPRESS RESULTS ACROSS VARIOUS SAMPLES (SEQUEL II SYSTEM)

Sample #	Sample Description	Protocol	# FLNC Reads	% FLNC Reads
1	UHR	Standard	3,466,513	85%
2	Mouse Liver	Standard	3,431,638	87%
3	MCF7	Standard	3,531,419	84%
4	Brain	Standard	2,943,148	86%
5	Alz Brain Tissue	Standard	3,142,634	83%
6	Heart	Standard	2,753,509	87%
7	Liver	long	3,542,983	85%
8	ColT Cell Line	short	2,852,434	84%

ארק כל ארכן כ



- Iso-Seq Universal Human Control

Metric	
Number of Raw Bases (Gb)	255
Total Reads	4,460,955
Full Length Non-chimeric Reads	3,436,022
CCS Passes (Mean)	8



Data shown above from a Universal Human Reference RNA (human) and Lexogen SIRV spike-in controls. The library was constructed using the SMRTbell Template Prep Kit on a Sequel II System (1.0 Chemistry, Sequel II System Software v7.0, 20-hour movie). Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.



Sequencing for High-Resolution Metagenomics and Closed Microbial Genomes סיק כל ארכן כל

#### PACBIO PROVIDES HIGHER-RESOLUTION VIEWS OF BOTH METAGENOME MEMBERS AND FUNCTIONS

#### Full-length 16S provide species and strain-level information

- Highly accurate, full-length 16S sequences distinguishes keystone or critical species from genus-level noise
- Follow our <u>16S amplification and sequencing protocol</u> or use an all-in-one kit from our partner <u>Shoreline Biome</u>

# HiFi Shotgun profiling reveals intact genes and operons without assembly

- Move beyond correlation and reveal community functions with highly accurate, 10 kb HiFi reads
- Obtain high-confidence information even for low-abundance sequences seen only once

#### Metagenome assembly from long reads generates new references

- Resolve genomes of microbes that can't be easily cultured
- Leverage epigenomic data to cluster contigs and plasmids from the same strain

סאכן כל אכן כל איכ

# ONLY PACBIO CAN DELIVER SINGLE MOLECULE LONG READS WITH HIGH ACCURACY



- With 100 kb average read lengths, both 16S and shotgun libraries can be sequenced at very high accuracy on the Sequel II system
- 16S: Up to 2 M Q30 full-length 16S sequences
- Shotgun: Up to 1.7 M Q20 reads from 10 kb libraries



**HiFi Read Accuracy** 

#### **FULL-LENGTH 16S SEQUENCING**



Michael K Dougherty

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

This article is a preprint and has not been peer-reviewed [what does this mean?].

- Nearly every bacteria has multiple copies of the 16S housekeeping gene, but in many cases they are not perfect duplicates
- PacBio CCS produced multiple distinct 16S sequences per bacterial genome, and they appear in integer ratios that reflected their copy number in each genome

"The high resolution and accuracy we are reporting derives in part from the exceptional and not-entirelyappreciated accuracy of PacBio CCS sequencing."



Callahan, BJ et. al. (2018) High-throughput amplicon sequencing of the full-length 16S rRNA gene with single-nucleotide resolution. bioRxiv doi: <u>http://dx.doi.org/10.1101/39233</u>.

### סאק כוא כן כוא כן כוא כן כוא כ

# FULL-LENGTH 16S TYPICAL PERFORMANCE ON THE SEQUEL II SYSTEM

16S	>Q20 BC reads	>Q20 Read Quality	BC Samples / SMRT Cell	Avg reads / BC
MSA-1002	2,426,218	Q42	96	35,000
MSA-1003 (1a)	1,743,260	Q34	48	36,317
MSA-1003 (2a)	1,738,543	Q34	96	18,109

Full-length 16S sequencing and shotgun profiling data were collected for mock communities

- -ATCC 20 strain staggered (MSA-1002)
- -ATCC 20 strain even (MSA-1003)

PACBIO FULL-LENGTH 16S SEQUENCING FAITHFULLY RECAPITULATES A MOCK COMMUNITY SAMPLE



Rhodobacter sphaeroides (ATCC 17029) Deinococcus radiodurans (ATCC BAA-816) Pseudomonas aeruginosa (ATCC 9027) Actinomyces odontolyticus (ATCC 17982) Bifidobacterium adolescentis (ATCC 15703) Cutibacterium acnes (ATCC 11828) Neisseria meningitidis (ATCC BAA-335) Escherichia coli (ATCC 700926) Porphyromonas gingivalis (ATCC 33277) Bacteroides vulgatus (ATCC 8482) Acinetobacter baumannii (ATCC 17978) Enterococcus faecalis (ATCC 47077) Streptococcus agalactiae (ATCC BAA-611) Streptococcus mutans (ATCC 700610) Bacillus cereus (ATCC 10987) Lactobacillus gasseri (ATCC 33323) Staphylococcus aureus (ATCC BAA-1556) Helicobacter pylori (ATCC 700392) Staphylococcus epidermidis (ATCC 12228)

PAC**BIO**®

Clostridium beijerinckii (ATCC 35702)

V1-V9 amplicons were sequenced on a single SMRT Cell 8M at 96-plex

ארק כל אכן כל איכן כל איכן כל איכן כל איכן כל א

# 16S: UNBIASED TAXONOMIC REPRESENTATION COMPARING TO ZYMOBIOMICS MICROBIAL COMMUNITY STANDARD



- Full-length 16S sequencing on the Sequel II System allows you to more faithfully characterize the composition of your metagenome
- With our internal protocol, 48-plex runs reliably generate at least 10 K reads / barcode. At 96-plex, ~5% of barcodes drop below this floor.



### סיק כל יכן כל יכן כל יכן איר פין כל ייכן איר פין אראי פין איר פין פאראיניי א 🔊 PAC**BIO**\*

#### LONG READS + HIGH ACCURACY MEANS GENE DISCOVERY CAN BE DONE *DIRECTLY* ON HIFI READS, WITHOUT ASSEMBLY

Sample	# HiFi Reads	# Amino acid sequences	Genes / read	Mean protein size (amino acids)
Human fecal 1	2,802,471	16,429,903	5.9	320.2
Human fecal 2	1,646,208	11,993,089	7.3	325.1
Human fecal 3	1,593,641	13,054,811	8.2	321.7

- FragGeneScan was used to predict proteins directly on unassembled HiFi reads
- Error-free genes can be found even from species with too little coverage for assembly
- High accuracy means existing NGS tools and pipelines can be used without modification

סאק כלא כן כל איכן כל איכן כל איכן כל איכן כל איכ

# PACBIO IS THE NEW GOLD STANDARD FOR BACTERIAL GENOMICS



Growth in published microbial genomes completed with SMRT Sequencing



Tanizawa et al. (2015) <u>Complete genome sequence and analysis of Lactobacillus hokkaidonensis LOOC260T, a psychrotrophic lactic acid bacterium</u> isolated from silage. BMC Genomics 16: 240.

# AFFORDABLY CHARACTERIZE COMPLETE MICROBIAL GENOMES ON THE SEQUEL PLATFORM



 Multiplex up to 30 Mb of microbial genomes / 16 samples on one SMRT Cell 1M

PACBIO<sup>®</sup>

- Adjust planned multiplexing depth to balance cost constraints with your requirements for genome completeness
- Use our Microbial Multiplexing Calculator to simplify equimolar pooling
- Assemble most bacterial chromosomes into 5 contigs or fewer

### סאק כלא כין כלא כין כלא כין כלא כין כלא כין כלא איז סיין כלא פא**פו**ס" PAC**BIO**"

#### ACHIEVE HIGHER COST EFFICIENCY AFTER BUILDING INITIAL EXPERIENCE WITH 30 MB POOLED LIBRARIES

Barcode ID	Sample ID	Pre-Assembly Yield (%)	GC Content (%)	Genome Size (bp)	Contigs (#)	Concordance w. Reference (QV)
BC1002	E. coli control 1	77	50.08	4,642,523	1	55.21
BC1015	E. coli control 2	77.9	50.08	4,642,523	1	56.25
BC1004	B. sub control 1	76.4	43.94	4,045,593	1	51.16
BC1016	B. sub control 2	74.9	43.94	4,045,593	1	51.45
BC1009	E. coli 1	77.4	50.08	4,642,523	1	54.63
BC1018	E. coli 2	77.7	50.08	4,642,523	1	54.91
BC1012	S. sonnei 1	76.4	51.03	4,813,450	1	59.83
BC1020	S. sonnei 2	78.1	51.03	4,813,450	1	54.27
BC1014	L. monocytogenes 1	71.8	37.94	3,032,269	1	53.06
BC1022	L. monocytogenes 2	73.2	37.94	3,032,269	1	49.63

- 10-plex run of 42.4 MB total microbial genomes shown

- Combined use of Barcoded Adapter Kit 8A & Barcoded Adapter Kit 8B
- Less than 5 contig assemblies for main chromosomal genomes can be achieved with Advanced HGAP parameters
- QV50 is roughly 99.999% base accuracy



## Sequel II System Updates

סאק כלא כן כלא כן כלא כן כלא כן כלא כ

#### SEQUEL II CHEMISTRY 2.0 (EARLY ACCESS (EA) INITIATED END OF JUNE)

- Higher yield of long HiFi reads due to greater efficiency of getting reads into rolling circle synthesis and continuing to the end of the movie
- Translates to ~25-50% average RL increase and ~40% higher yield of reads > Q20 and Q30
- Also allows increase in insert size from ~11-13 to 15-20+ kb, depending on application



Site	Library Name	Sample Biology	Insert Size (Kb)	[Loading]	Yield (Gb)	Pol RL	Subread	PO	P1	Base Rate	≥Q20 Yield (Gb)	≥ Q20 Reads (count)	≥ Q20 Read Length (mean, bp)	≥ Q20 Read Quality (median)	Percent Reads ≥ Q20	UMY (Gb)
Site 3 *	GM19734	Human	16.1	50pM	384.9	98259	11162	49%	49%	1.97	51. <b>2</b>	4,643,311	11029	Q36		
	GM19734	Human	16.1	60pM	437.3	96301	11595	41%	57%	1.95						
	GM19030	Human	22.8	60pM	380.8	93318	10222	47%	51%	1.97	47.1 4620262	10170	026			
	GM19030	Human	22.8	60pM	414.7	91150	10376	40%	57%	1.94	47.1	4029302	10170	0,50		
	HG03736	Human	19.1	60pM	388.3	87645	12765	42%	55%	1.90	52.2	2 016 221	12227	022		
	HG03736	Human	19.1	60pM	367.8	85863	12833	44%	53%	1.86	52.2	3,910,231	15332	U33		

Half of bases in reads: >190 kb



2,500,000



PAC**BIO**®



Data shown above from a 20 kb size-selected human library using the SMRTbell Template Prep Kit on a Sequel II System (2.0 Chemistry, Sequel II System Software v8.0, 30-hour movie). Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.
ארק כל ארכן כל ארכין כל ארכי

### HIFI SEQUENCING PERFORMANCE

# Q20 (99%) single-molecule accuracy



Read Quality Score (Q)

Data shown above from a 20 kb size-selected human library using the SMRTbell Template Prep Kit on a Sequel II System (2.0 Chemistry, Sequel II System Software v8.0, 30-hour movie). Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.

וליכן כלי כן כל יכן יכל יכ

PAC**BIO**®



Data shown above from a 15 kb size-selected human library using the SMRTbell Template Prep Kit on a Sequel II System (2.0 Chemistry, Sequel II System Software v8.0, 30-hour movie). Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.

ארק כל ארכן כל ארכין כל ארכי

### HIFI SEQUENCING PERFORMANCE



Data shown above from a 15 kb size-selected human library using the SMRTbell Template Prep Kit on a Sequel II System (2.0 Chemistry, Sequel II System Software v8.0, 30-hour movie). Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.

USING HIFI READS FOR VARIANT DISCOVERY Chemistry 1.0 Chemistry 2.0



Variant calls from ~15-fold HiFi read coverage of a human genome (HG002) were measured against the Genome in a Bottle small variant benchmark (v3.3.2) for SNVs and indels using Deep Variant. Libraries were generated using an 11 kb insert for Chemistry 1.0 and 15 kb insert for Chemistry 2.0.

**DE NOVO ASSEMBLY WITH HIFI READS – RICE GENOME** 

Library Size	17 kb	24 kb
Raw Yield	218 Gb	366 Gb
1 SMRT Cell Yield	15 Gb	25 Gb
Median Read Length	16 kb	21 kb
Coverage (400Mb)	38-fold	63-fold
FALCON Asm Length	403 Mb	405 Mb
N Contigs	209	211
Contig N50	14 Mb	20 Mb
N Chrom in 1 Contig	0	3





PAC**BIO**°

Data shown above from a 17 kb & 24 kb size-selected rice library using the SMRTbell Template Prep Kit on a Sequel II System (2.0 Chemistry, Sequel II System Software v8.0, 30-hour movie). Falcon assembly was performed post CCS analysis. Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.





Data shown above from a Universal Human Reference RNA (human) and Lexogen SIRV spike-in controls. The library was constructed using the Iso-Seq Express workflow including the SMRTbell Express Template Prep Kit 2.0 on a Sequel II System (Sequel II Sequencing kit 2.0, Sequel II Binding Kit 2.0, Sequel II System Software v8.0, 24-hour movie). Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.

ליכר כד כד כד כ

PACBIO<sup>®</sup>



Data shown above from a 35 kb size-selected *E. coli* library using the SMRTbell Template Prep Kit on a Sequel II System (2.0 Chemistry, Sequel II System Software v8.0, 15-hour movie). Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.



## SMRT Link v8.0 Preview

### ארק כל אכן כל איכ

#### SMRT LINK V8.0 – SUMMARY OVERVIEW OF MAJOR IMPROVEMENTS

#### **New Workflow Engine - Cromwell**

- Scalable, widely adopted in scientific environments, developed by Broad

#### **Further Enhancements for CCS Analysis**

- Faster CCS analysis (4-6 times faster time-to-results)

#### **Analysis Applications**

- NEW: Microbial Assembly application
- de novo assembly for CCS data Falcon Unzip (Bioconda)

#### **UI and Usability improvements**

- New style GUI same functionality
- Support for IGV visualization
- Analysis restart from the point of failure





#### **ANALYSIS APPLICATIONS**

X×x

#### **Microbial Assembly – NEW application**

- de novo assembly for small genomes between 1.9-10 Mb
- Circularization
- Plasmid assembly (2 220 kb)



#### Iso-Seq Analysis

- Analysis from CCS data



#### de novo assembly with HiFi reads - Bioconda

- Direct CCS (FASTA+FASTQ) input to FALCON and FALCON-Unzip
- Accurate phasing good phased contig QV approaching Q50
- Faster read tracking and polishing using CCS reads

## סאכן כל אכן כל איכן כל איכן כל איכן כל איכן כל א

#### **COMPATIBILITY AND SUPPORT**

#### SMRT Link v8.0 supports Sequel and Sequel II Systems only

- Use SMRT Link v7.0 or earlier for PacBio RS II data

#### **Data Visualization**

- IGV compatibility URL links to alignment results (BAM and BAI)
- SMRT View is no longer included with SMRT Link

#### SMRT Link v8.0 migration

- This is a major release more efforts for customers to upgrade
  - Relatively minimal impact on UI users, cmd and APIs users more impacted
- Sept.11, 2019 SMRT Link v8.0 notification email details on upcoming changes

ארק כל אכן כל איכ

### WHAT COULD YOU DO WITH A SEQUEL II SYSTEM TODAY?

#### With a single SMRT Cell 8M:

- Generate a 2 Gb genome assembly
- Call structural variants in a human genome
- Sequence a whole transcriptome
- Determine the composition of a >90 gut microbiome samples



#### With 2-3 SMRT Cells 8M:

- Detect all variants in a human genome
- Phase a diploid assembly of a human genome



Plant and vertebrate genome assemblies generated with SMRT Sequencing data from the genome database on NCBI showing PacBio assemblies readily achieve contig N50s ≥1 Mb.

Coverage of the Cytochrome P450 2D6 (CYP2D6) and CYP2D7 genes with HiFi reads and NGS reads visualized in IGV. CYP2D6 is responsible for the metabolism and elimination of approximately 25% of clinically used drugs.



#### 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. Pacific Biosciences does not sell a kit for carrying out the overall No-Amp Targeted Sequencing method. Use of this No-Amp method may require rights to third-party owned intellectual property. BluePippin and SageELF are trademarks of Sage Science. NGS-go and NGSengine are trademarks of GenDx. FEMTO Pulse and Fragment Analyzer are trademarks of Agilent Technologies Inc.

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