

## Comprehensive Views of Genetic Diversity with Single Molecule, Real-Time (SMRT) Sequencing

#### Alix Kieu Cruse

November 2015

For Research Use Only. Not for use in diagnostics procedures. © Copyright 2015 by Pacific Biosciences of California, Inc. All rights reserved.

## AGENDA

SMRT Technology Overview
High quality reference genomes

Microbiology and Infectious Diseases
Animal and Plant Genomes
Structural Variation detection

Transcriptomic approaches – IsoSeq
Epigenomics – base modification detection

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

о и со со со со со со со со со расвю<sup>•</sup>

#### Zero-Mode Waveguides





**Phospholinked** 

Up to 1 million ZMWs per SMRT<sup>®</sup> Cell



SMRT Cells containing up to a million ZMWs are processed on PacBio<sup>®</sup> Systems which simultaneously monitor each of the waveguides in real time.

KEY SEQUENCING CHARACTERISTICS

ליכן כוי כו יכן כו יכן כו יכן כו יכ

#### **1.** Contiguity

- Sequence reads >10,000 bases
- Some reads >50 kb

#### **2.** Accuracy

- Achieves >99.999% (QV50)
- Lack of systematic sequencing errors

#### **3.** Uniformity

- Lack of GC content or sequence

complexity bias

#### **4.** Originality

- No DNA amplification
- Epigenome characterization



PACBIO

### ס-רן כל-רכן כל-רכן כל -רכן כל -

#### **CONSENSUS ACCURACY PERFORMANCE COMPARISON**



~QV 50 (99.999%) consensus accuracy coverage:

- P6-C4: 30x ±10
- P5-C3: 45x ±10

Reduced coverage

requirements & throughput

increases combined result in

overall 3-4x performance

#### improvement

E. coli 20kb-insert library, resequencing analysis with SMRT Analysis v2.3

PCR-FREE SAMPLE PREP WORKFLOW MEANS LESS BIAS

ס- כן כן כן כן כן כן כן כן כן ארכין כן ארכין כן ארכין כן ארכין כן ארכין כ

#### **Sample Preparation**

Building of the SMRTbell template



### SAMPLING OF APPLICATION REQUIREMENTS

	Input Material	Sample Processing	
	50ng – 1ug	One 10 kb Library prep 1-2 SMRT <sup>®</sup> Cells / 5 MB	
Bacteria		SMRT HGAP Portal Analysis	
Iso-Seq	Total RNA>>cDNA 5ng-1ug	Project Scope Specific: Whole Transcriptome Discovery 2-3 Libraries Size Selection with Blue Pippin SMRT Portal IsoSeq Analysis	
Long Amplicons (HLA)	50-500ng	1 library Capabilities for multiplexing SMRT Portal Long Amplicon Analysis Third Party – GenDX /Conexio	All use standard SMRTbe prep
Lower Eukaryotes	5-10ug	One 20 kb Library prep Size selection with Blue Pippin SMRT Portal HGAP Analysis	
-	20-100ug	One to three 20 kb Library prep Size selection with Blue Pippin SMRT Portal FALCON	

Higher Eukaryotes

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

#### SMRT SEQUENCING: BROADLY ACCEPTED AND USED

- **150+** PacBio sequencing sites worldwide
- **900+** Publications to date with PacBio data
- **1-2** New publications average *per day*

ס-רק כל ארכן כל

#### CONTIG N50 ASSEMBLY STATISTICS OF GENOMES ASSEMBLED USING ONLY PACBIO DATA



# MICROBIOLOGY AND INFECTIOUS DISEASES APPLICATIONS

PACBIO®

 $\bigcirc$ 

#### סארק כל ארכן כל ארכין כל ארכי

#### WHY GENOME FINISHING?

#### Clostridium autoethanogenum

Industrial microbe for commodity chemicals

- PacBio-only assembly closed, high-quality genome without manual finishing
- Full metabolic pathway reconstruction
- Complete genome revealed genes important for biofuel production missed by short read technologies (~100 contigs)



Brown et al. (2014) Comparison of single-molecule sequencing and hybrid approaches for finishing the genome of Clostridium autoethanogenum and analysis of CRISPR systems in industrial relevant Clostridia. *Biotechnology for Biofuels* **7**:40

#### RESEARCH ARTICLE

#### ANTIBIOTIC RESISTANCE

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

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

Public health officials have raised concerns that plasmid transfer between Enterobacteriaceae species may spread resistance to carbapenems, an antibiotic dass of last resort, thereby rendering common health care-associated infections nearly impossible to treat. To determine the diversity of carbapenemase-encoding plasmids and assess their mobility among bacterial species, we performed comprehensive surveillance and genomic sequencing of carbapenemresistant Enterobacteriaceae in the National Institutes of Health (NIH) Clinical Center patient population and hospital environment. We isolated a repertoire of carbapenemase-encoding Enterobacteriaceae, including multiple strains of Klebsiella pneumoniae, Klebsiella oxytoca, Escherichia coli, Enterobacter cloacae, Citrobacter freundii, and Pantoea species. Long-read genome sequencing with full end-to-end assembly revealed that these organisms carry the carbapenem resistance genes on a wide array of plasmids. K. pneumoniae and E. cloacae isolated simultaneously from a single patient harbored two different carbapenemase-encoding plasmids, indicating that plasmid transfer between organisms was unlikely within this patient. We did, however, find evidence of horizontal transfer of carbapenemaseencoding plasmids between K. pneumoniae, E. doacae, and C. freundii in the hospital environment. Our data, including full plasmid identification, challenge assumptions about horizontal gene transfer events within patients and identify possible connections between patients and the hospital environment. In addition, we identified a new carbapenemaseencoding plasmid of potentially high clinical impact carried by K. pneumoniae, E. coli, E. doacae, and Pantoea species, in unrelated patients and in the hospital environment.



Presented at the 60<sup>th</sup> NIH Clinical Center Anniversary, November 6<sup>th</sup>, 2013



National Human Genome Research Institute

# Patients and hospital environmental isolates: horizontal gene transfer of carbapemase resistance encoding plasmid?



# Gallery of Plasmids



31 plasmids (27 unique) across the 10 isolates sequenced

KPC-2KPC-3

# Full genome sequencing rules out patient-patient transmission



Targeted sequencing could have given misleading result



# Horizontal gene transfer from patient to environment







# PLANT AND ANIMAL GENOMES



ס-רק כל ארכן כל

#### CONTIG N50 ASSEMBLY STATISTICS OF GENOMES ASSEMBLED USING ONLY PACBIO DATA



סאכן כל אכן כל אכן כל אכן כל איכן כל אי

#### **OROPETIUM GENOME FOR DROUGHT STUDY**



### **Robert VanBuren & Todd Mockler**

Donald Danforth Plant Science Center, St. Louis





#### Oropetium "resurrection plant", 250 Mb genome

#### **COMPARISON OF ORO GENOME ASSEMBLIES**

#### Illumina

#### 6 Small Insert Libraries:

180bp insert PE 2X100bp
250bp insert PE 2X100bp
500bp insert PE 2X100bp
500bp insert PE 2X250bp
1kb insert PE 2X100bp
1kb insert PE 2x300bp

#### **4 Mate Pair Libraries**

3kb insert PE 2X76bp 8kb insert PE 2X76bp 9kb insert PE 2X76bp 18kb insert PE 2X76bp

#### **PacBio**

#### **One 20 KB Insert Library:**

BluePippin<sup>™</sup> Size Selection @ 15 kb Protocol P6-C4 Chemistry Mean read length : 12,872 bp N50 read length : 16,485 10x genome coverage >20kb

Illumina <sup>®</sup> Assembly	y Statistics
--------------------------------	--------------

# Scaffolds	14,216
Scaffold N50	11 kb
Total size	158 Mb
Total # Ns	1.4Mb

st genome	assembled	63%

PacBio®	Assembly	<b>Statistics</b>	(HGAP)
---------	----------	-------------------	--------

Fet genome assembled	47 %
Danaat Contant	17 %
Sum of Contig Lengths	244.46 Mb
Max Contig Length	7.98 Mb
Contig N50	2.38 Mb
# Polished Contigs	625

### סיכן כל יכן כל יכן כל יכן כל יכן כל יכן כל יכ

Protein coding regions

Gap regions

#### COMPARATIVE ANALYSIS HIGHLIGHTS SYNTENY AND GAP REGIONS

#### Oropetium Genome



#### **RESOLVING TANDEM REPEATS IN DROSOPHILA Y CHR**

G3: Genes|Genomes|Genetics Early Online, published on April 9, 2015 as doi:10.1534/g3.115.017277

accepted v. 6apr2015 10AM

Long-read single molecule sequencing to resolve tandem gene copies: The *Mst*77Y region on the *Drosophila melanogaster* Y chromosome

Flavia J. Krsticevic<sup>\*</sup>, Carlos G. Schrago<sup>§</sup> and A. Bernardo Carvalho<sup>§</sup>



ס- רן כל - רן כל - רן כל - ר כן כל - ר כ<mark>ן כל - ר</mark>פי כל - ר

#### ASSEMBLY COMPARISON OF MST77Y GENES

TABLE 1. Mst77Y genes in different assemblies of the D. melanogaster genome.

Assembly	Mst77Y	Perfect	With errors	Number of	Scaffold size (kb)
	genes found	matches <sup>a</sup>		scaffolds	
SLR	10	8	2	7	3 - 13
MHAP	18	18	-	1	747
PBcR	20	17	3	2	20 ; 177
FALCON	18	11	7	1	619
WGS3	6	2	4	6	< 2



Krsticevic, F. et al. G3. April 9, 2015, doi: 10.1534/g3.115.017277

#### ASSEMBLY IMPROVEMENTS WITH PBJELLY 2

PBJelly Examples	D. pseudoobscura Fly	<i>S. purpuratus</i> Sea urchin	<i>L. variegatus</i> Sea urchin	<i>P. anubis</i> Baboon	<i>R. norvegicus</i> Rat	C. <i>atys</i> Mangabey	<i>M. murinus</i> Mouse Lemur	<i>O.Aries</i> Sheep
Genome Size	146 Mb	0.8 Gb	0.8 Gb	2.8 Gb	3.0 Gb	2.8 Gb	2.8 Gb	2.8 Gb
PacBio Coverage	24x	10.6x	10.0x	11.5x	9.3x	12x	23x	19x
Closed Gaps	70%	38%	26%	64%	40%	63%	57%	89%
Starting Gaps	6,010	142,764	433,764	49,376	109,263	190,087	123,464	117,293
Remaining Gaps	1,793	87,844	321,392	17,667	65,639	70,875	52,634	12,766
Contig N50 Before	53 kb	l4 kb	6.3 kb	134 kb	60 kb	35 kb	56 kb	41.7 kb
Contig N50 After	215 kb	24 kb	14.5 kb	376 kb	131 kb	137 kb	191 kb	501 kb
	100	C. C	3-11 M 1 4		Z	141 200	A COLOR	314 ASS

- Improved contiguity
  - Closed up to 89% of gaps
  - Increased Contig N50 to 500 kb
- Improved transcript alignments



Kim C. Worley, Ph.D. Plant and Animal Genome XXII San Diego, CA, January 11-15, 2014



Lower quality draft genomes that start with a low contig N50 before addition of any PacBio<sup>®</sup> data show less improvement



# HUMAN BIOMEDICAL APPLICATIONS



#### STRUCTURAL VARIATION IS THE PREDOMINANT FORM OF SEQUENCE VARIATION IN THE HUMAN GENOME

OPEN aCCESS Freely available online

PLOS BIOLOGY

# The Diploid Genome Sequence of an Individual Human

Samuel Levy<sup>1\*</sup>, Granger Sutton<sup>1</sup>, Pauline C. Ng<sup>1</sup>, Lars Feuk<sup>2</sup>, Aaron L. Halpern<sup>1</sup>, Brian P. Walenz<sup>1</sup>, Nelson Axelrod<sup>1</sup>, Jiaqi Huang<sup>1</sup>, Ewen F. Kirkness<sup>1</sup>, Gennady Denisov<sup>1</sup>, Yuan Lin<sup>1</sup>, Jeffrey R. MacDonald<sup>2</sup>, Andy Wing Chun Pang<sup>2</sup>, Mary Shago<sup>2</sup>, Timothy B. Stockwell<sup>1</sup>, Alexia Tsiamouri<sup>1</sup>, Vineet Bafna<sup>3</sup>, Vikas Bansal<sup>3</sup>, Saul A. Kravitz<sup>1</sup>, Dana A. Busam<sup>1</sup>, Karen Y. Beeson<sup>1</sup>, Tina C. McIntosh<sup>1</sup>, Karin A. Remington<sup>1</sup>, Josep F. Abril<sup>4</sup>, John Gill<sup>1</sup>, Jon Borman<sup>1</sup>, Yu-Hui Rogers<sup>1</sup>, Marvin E. Frazier<sup>1</sup>, Stephen W. Scherer<sup>2</sup>, Robert L. Strausberg<sup>1</sup>, J. Craig Venter<sup>1</sup>

1 J. Craig Venter Institute, Rockville, Maryland, United States of America, 2 Program in Genetics and Genomic Biology, The Hospital for Sick Children, and Molecular and Medical Genetics, University of Toronto, Toronto, Ontario, Canada, 3 Department of Computer Science and Engineering, University of California San Diego, La Jolla, California, United States of America, 4 Genetics Department, Facultat de Biologia, Universitat de Barcelona, Barcelona, Catalonia, Spain

Presented here is a genome sequence of an individual human. It was produced from ~32 million random DNA fragments, sequenced by Sanger dideoxy technology and assembled into 4,528 scaffolds, comprising 2,810 million

Sec

fro

seq

Seq

from

bases (Mb) of contiguous sequence with approximately 7.5-fold commodified version of the Celera assembler to facilitate the identification individual diploid genome. Comparison of this genome and the *b* reference assembly revealed more than 4.1 million DNA varian 1,288,319 were novel) included 3,213,401 single nucleotide poly bp), 292,102 heterozygous insertion/deletion events (indels)(1-inversions, as well as numerous segmental duplications and coarcounts for 22% of all events identified in the donor, however important role for non-SNP genetic alterations in defining the di heterozygous for one or more variants. Using a novel haplotyr genome sequence in segments >200 kb, providing further prec depict a definitive molecular portrait of a diploid human geno comparisons and enables an era of individualized genomic info

Citation: Levy S, Sutton G, Ng PC, Feuk L, Halpern AL, et al. (2007) The diploid genome seq 0050254

#### Introduction

Each of our genomes is typically composed of DNA packaged into two sets of 23 chromosomes; one set inherited from each parent whose own DNA is a mosaic of preceding ancestors. As such, the human genome functions as a diploid entity with phenotypes arising due to the sometimes complex interplay of alleles of genes and/or their noncoding functional regulatory elements.

The diploid nature of the human genome was first observed as unbanded and banded chromosomes over 40 years ago [1– 4], and karyotyping still predominates in clinical laboratories as the standard for global genome interrogation. With the "Non-SNP DNA variation accounts for 22% of all events identified in the donor, however they involve 74% of all variant bases. This suggests an important role for non-SNP genetic alterations in defining the diploid genome structure."

PACBIO

version of the genome is a consensus sequence derived from five individuals. Both versions almost exclusively report DNA variation in the form of single nucleotide polymorphisms (SNPs). However smaller-scale (<100 bp) insertion/deletion sequences (indels) or large-scale structural variants [10–15] also contribute to human biology and disease [16–18] and warrant an extensive survey.

#### PACBIO DEFINES THE NEW GOLD STANDARD IN HUMAN GENOME DE NOVO ASSEMBLY

PACBIO



Data sources: HuRef (Venter) (http://www.plosbiology.org/article/info:doi/10.1371/journal.pbio.0050254); BGI YH (http://genome.cshlp.org/content/20/2/265.abstract Table II); KB1 (http://www.nature.com/nature/journal/v463/n7283/full/nature08795.html); NA12878 (http://www.pnas.org/content/early/2010/12/20/1017351108.abstract Table3); CHM1 Illumina (http://www.ncbi.nlm.nih.gov/assembly/GCF\_000306695.2/)

#### HUMAN GENOME SEQUENCING WITH SMRT TECHNOLOGY

NATURE | LETTER

doi:10.1038/nature13907

PACBIO

# Resolving the complexity of the human genome using single-molecule sequencing

Mark J. P. Chaisson<sup>1</sup>, John Huddleston<sup>1,2</sup>, Megan Y. Dennis<sup>1</sup>, Peter H. Sudmant<sup>1</sup>, Maika Malig<sup>1</sup>, Fereydoun Hormozdiari<sup>1</sup>, Francesca Antonacci<sup>3</sup>, Urvashi Surti<sup>4</sup>, Richard Sandstrom<sup>1</sup>, Matthew Boitano<sup>5</sup>, Jane M. Landolin<sup>5</sup>, John A. Stamatoyannopoulos<sup>1</sup>, Michael W. Hunkapiller<sup>5</sup>, Jonas Korlach<sup>5</sup> & Evan E. Eichler<sup>1,2</sup>

The human genome is arguably the most complete mammalian reference assembly<sup>1-3</sup>, yet more than 160 euchromatic gaps remain<sup>4-6</sup> and aspects of its structural variation remain poorly understood ten years after its completion7-9. To identify missing sequence and genetic variation, here we sequence and analyse a haploid human genome (CHM1) using single-molecule, real-time DNA sequencing<sup>10</sup>. We close or extend 55% of the remaining interstitial gaps in the human GRCh37 reference genome-78% of which carried long runs of degenerate short tandem repeats, often several kilobases in length, embedded within (G+C)-rich genomic regions. We resolve the complete sequence of 26,079 euchromatic structural variants at the base-pair level, including inversions, complex insertions and long tracts of tandem repeats. Most have not been previously reported, with the greatest increases in sensitivity occurring for events less than 5 kilobases in size. Compared to the human reference, we find a significant insertional bias (3:1) in regions corresponding to complex insertions and long short tandem repeats. Our results suggest a greater complexity of the human genome in the form of variation of longer and more complex repetitive DNA that can now be largely resolved with the application of this longer-read sequencing technology.

LETTEI

a template for recruiting additional sequence reads for assembly (Supplementary Information). Using this approach, we closed 50 gaps and extended into 40 others (60 boundaries), adding 398 kb and 721 kb of novel sequence to the genome, respectively (Supplementary Table 4). The closed gaps in the human genome were enriched for simple repeats, long tandem repeats, and high (G+C) content (Fig. 1), but also included novel exons (Supplementary Table 20) and putative regulatory sequences based on DNase I hypersensitivity and chromatin immunoprecipitation followed by high-throughput DNA sequencing (ChIP-seq) analysis (Supplementary Information). We identified a significant 15-fold enrichment of short tandem repeats (STRs) when compared to a random sample (*P* < 0.00001) (Fig. 1a). A total of 78% (39 out of 50) of the closed gap sequences were composed of 10% or more of STRs. The STRs were frequently embedded in longer, more complex, tandem arrays of degenerate repeats reaching up to 8,000 bp in length (Extended Data Fig. 1a-c), some of which bore resemblance to sequences known to be toxic to Escherichia coli<sup>16</sup>. Because most human reference sequences<sup>17,18</sup> have been derived from clones propagated in E. coli, it is perhaps not surprising that the application of a long-read sequence technology to uncloned DNA would resolve such gaps. Moreover, the length and

Chaisson et al. (2014) Nature doi:10.1038/nature13907

סאכן כל אכן כל איכ

#### PACBIO DATA VS. GRCH37 & 1000 GENOMES PROJECT

- Closed 55% of interstitial gaps remaining in reference genome
- Resolved 26,079 euchromatic structural variants at the base-pair level
- ~22,000 (85%) of these are novel
- 6,796 of the events map within 3,418 genes



Chaisson et al. (2014) Nature doi:10.1038/nature13907

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

#### PACBIO DATA VS. PRESENT-DAY ILLUMINA DATA

# LETTER

doi:10.1038/nature13907

# Resolving the complexity of the human genome using single-molecule sequencing

Mark J. P. Chaisson<sup>1</sup>, John Huddleston<sup>1,2</sup>, Megan Y. Dennis<sup>1</sup>, Peter H. Sudmant<sup>1</sup>, Maika Malig<sup>1</sup>, Fereydoun Hormozdiari<sup>1</sup>, Francesca Antonacci<sup>3</sup>, Urvashi Surti<sup>4</sup>, Richard Sandstrom<sup>1</sup>, Matthew Boitano<sup>5</sup>, Jane M. Landolin<sup>5</sup>, John A. Stamatoyannopoulos<sup>1</sup>, Michael W. Hunkapiller<sup>5</sup>, Jonas Korlach<sup>5</sup> & Evan E. Eichler<sup>1,2</sup>

"Notably, less than 1% of these variant are present in newer assemblies of the human genome, including GRCh38 and CHM1.1 (ref. 22) (derived primarily by Illumina sequencing technology)."

Supplementary Table 7. Presence of inserted sequences in other human assemblies.

	Insert	ion	hg19		GRCh	38	CHM1	.1
	count	bases	count	bases	count	bases	count	bases
Complex	1116	2148286	0	0	4	21567	11	36303
STR	406	826962	0	0	0	0	2	2344
VNTR	215	498362	1	1846	1	1846	1	1846

Chaisson et al. (2014) Nature doi:10.1038/nature13907

#### **GENOME-WIDE STRUCTURAL VARIATION CHARACTERIZATION**

איק כלי כן כלי כן כלי כן כלי כ

English et al. BMC Genomics (2015) 16:286 DOI 10.1186/s12864-015-1479-3



PACBIO

#### **RESEARCH ARTICLE**



# Assessing structural variation in a personal genome—towards a human reference diploid genome

Adam C English<sup>1+</sup>, William J Salerno<sup>1\*+</sup>, Oliver A Hampton<sup>1</sup>, Claudia Gonzaga-Jauregui<sup>2</sup>, Shruthi Ambreth<sup>1</sup>, Deborah I Ritter<sup>1</sup>, Christine R Beck<sup>2</sup>, Caleb F Davis<sup>1</sup>, Mahmoud Dahdouli<sup>1</sup>, Singer Ma<sup>3</sup>, Andrew Carroll<sup>3</sup>, Narayanan Veeraraghavan<sup>1</sup>, Jeremy Bruestle<sup>4</sup>, Becky Drees<sup>4</sup>, Alex Hastie<sup>5</sup>, Ernest T Lam<sup>5</sup>, Simon White<sup>1</sup>, Pamela Mishra<sup>1</sup>, Min Wang<sup>1</sup>, Yi Han<sup>1</sup>, Feng Zhang<sup>6</sup>, Pawel Stankiewicz<sup>2</sup>, David A Wheeler<sup>1,2</sup>, Jeffrey G Reid<sup>1</sup>, Donna M Muzny<sup>1,2</sup>, Jeffrey Rogers<sup>1,2</sup>, Aniko Sabo<sup>1,2</sup>, Kim C Worley<sup>1,2</sup>, James R Lupski<sup>1,2,7,8</sup>, Eric Boerwinkle<sup>1,9</sup> and Richard A Gibbs<sup>1,2</sup>

#### Abstract

**Background:** Characterizing large genomic variants is essential to expanding the research and clinical applications of genome sequencing. While multiple data types and methods are available to detect these structural variants (SVs), they remain less characterized than smaller variants because of SV diversity, complexity, and size. These challenges are

### ארק כל ארכן כל ארכין כל ארכי

#### **GENOME-WIDE STRUCTURAL VARIATION CHARACTERIZATION**

- Structural variation survey on diploid genome
- Integrates different sequencing methods as well as BioNano optical mapping
- Integrated analysis pipeline described, available in cloud-based DNAnexus environment

Data	Туре	Resolution	Source
WGS Illumina HiSeq	NGS	48X 100x100 bp paired-end	[26]
WGS Illumina Nextera	NGS	2X 100x100 bp 6.5 kbp mate-pair inserts	Methods
WGS SOLID	NGS	3X 35 bp fragment 10X 25x25 bp paired-end 17X 50x50 bp paired-end	[25,26]
WGS PacBio	Long-Read	10X~10,000 bp	Methods
Agilent 1 M	aCGH	1-million-probe oligo array	[26]
NimbleGen 2.1 M	aCGH	2.1-million-probe oligo array	[26]
NimbleGen 4.2 M	aCGH	4.2-million-probe oligo array	Methods
Custom Agilent Exon Array	aCGH	44,000 neuropathy-specific oligo array	[26]
BioNano Irys	Genome Mapping	Single-molecule genome architecture	Methods
anger-Validated Deletions	Manual	42 fully resolved deletions	Methods

"Here, we characterize the SV content of a personal genome with Parliament, a publicly available consensus SV-calling infrastructure that merges multiple data types and SV detection methods." ארק כל ארכין כל ארכי

#### **GENOME-WIDE STRUCTURAL VARIATION CHARACTERIZATION**

 Detecting structural variation with Illumina is difficult, even when integrating different paired-end (PE) data:

Table 4 Illumina-only method comparison							
Program	<b>Total Called</b>	Supported	Unsupported	FDR	Sensitivity		
CNVnator	6,197	1,211	4,986	80.46%	22.62%		
BreakDancer	5,520	2,269	3,251	58.89%	42.39%		
Delly	3,720	1,669	2,051	55.13%	31.18%		
Crest	2,219	1,889	330	14.87%	35.29%		
Pindel	4,451	3,035	1,416	31.81%	56.70%		
SV-STAT	892	876	16	1.79%	16.36%		
Tiresias	1,347	417	930	69.04%	7.79%		
Spiral	1,881	1,824	57	3.03%	34.07%		
Parliament	3,082	2,852	258	8.37%	57.34%		

Performance for each Illumina-only method is summarized. Supported and Unsupported columns indicate the number of calls with and without local hybrid assembly support, respectively. False discovery rate (FDR) and sensitivity are calculated using all 17,704 Illumina Only reference inconsistent loci and the subset of 5,584 that are supported by hybrid assembly.

"Despite these benefits of a multi-algorithm approach, Illuminaonly discovery still only recovers approximately half of the 9,777 SVs identified by multi-source Parliament: PBHoney alone identifies 4,268 SVs supported by hybrid assembly, representing events "invisible" to PE data." 

#### **GENOME-WIDE STRUCTURAL VARIATION CHARACTERIZATION**

- With only 10x PacBio coverage:



"Applying multiple Parliament workflows, we demonstrate that while method integration is optimal for SV detection in Illumina paired-end data, the addition of long-read data can more than triple the number of SVs detectable in a personal genome."

# Iso-Seq: Full transcript sequencing



#### אק כוי כן כוי כן כן יכן יכ PACBIO®

### **DETERMINATION OF TRANSCRIPT ISOFORMS**



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

splice

#### "GENE IDENTIFICATION, EVEN IN WELL-CHARACTERIZED HUMAN CELL LINES AND TISSUES, IS LIKELY FAR FROM COMPLETE"

לא ליק בוי כן בויי כן כן יכן יכ

#### PNAS

### Characterization of the human ESC transcriptome by hybrid sequencing

Kin Fai Au<sup>a</sup>, Vittorio Sebastiano<sup>0</sup>, Pegah Tootoonchi Afshar<sup>6</sup>, Jens Durruthy Durruthy<sup>b</sup>, Lawrence Lee<sup>d.</sup>e, Brian A. Williams<sup>†</sup>, Harm van Bake<sup>19</sup>, Eric E. Schadt<sup>9</sup>, Renee A. Reijo-Pera<sup>b</sup>, Jason G. Underwood<sup>d,h.</sup>1, and Wing Hung Wong<sup>6-1</sup>

<sup>1</sup>Department of Statistics and Department of Health Research and Policy, Stanford University, Stanford, C4 94305, <sup>1</sup>Center for Human Pluripotent Stem Cell Research and Education, Department of Obstetrias and Gynecology, Institute for Stem Cell Biology and Regenerative Medicine, Stanford University School of Medicine, Stanford, CA 94305; <sup>1</sup>Department of Electrical Engineering, School of Engineering, Stanford University, Stanford C4, 194305; <sup>4</sup>Pacific Biosciences of California, Meno Park, CA 94025; <sup>1</sup>Invite enc, San Francisco, CA 94107; <sup>1</sup>Division of Biology and Beckman Institute, California Institute of Technology, Pasadena, CA 91125; <sup>1</sup>Department of Genetics and Genomic Sciences, Mount Sinai School of Medicine, New York, NY 10029-6574; and <sup>1</sup>University of Washington, Department of Genome Sciences, Seattle, WA 98105

Contributed by Wing Hung Wong, November 5, 2013 (sent for review August 6, 2013)

Although transcriptional and posttranscriptional events are detected in RNA-Seq data from second-generation sequencing, fulllength mRNA isoforms are not captured. On the other hand, thirdgeneration sequencing, which yields much longer reads, has current limitations of lower raw accuracy and throughput. Here, we combine second-generation sequencing and third-generation sequencing with a custom-designed method for isoform identification and quantification to generate a high-confidence isoform dataset for human embryonic stem cells (hESCs). We report 8.084 RefSeq-annotated isoforms detected as full-length and an additional 5,459 isoforms predicted through statistical inference. Over one-third of these are novel isoforms, including 273 RNAs from gene loci that have not previously been identified. Further characterization of the novel loci indicates that a subset is expressed in pluripotent cells but not in diverse fetal and adult tissues; moreover, their reduced expression perturbs the network of pluripotency-associated genes. Results suggest that gene identification, even in well-characterized human cell lines and tissues, is likely far from complete

isoform discovery | PacBio | hESC transcriptome | alternative splicing | IncNRA

In the 5 y since the introduction of RNA-Seq (1, 2), there have been remarkable advances in our ability to analyze the transcriptome. During this period, additional methods based on nextgeneration sequencing (NGS) have been developed for the study of many different aspects of RNA biology. These include methods to study RNA species that are ribosome bound (3), nuclear (4), implicated in RNA editing (5), functional noncoding RNAs (6), in protein–RNA binding sites (7), and interacting in microRNA-mRNA complexes (8). Concurrently, the increase in NGS throughput and the development of multiplex sequencing protocols have made RNA-Seq analysis as cost-effective as gene expression microarrays.

Despite these advances, we are still far from achieving the original goals of RNA-Seq analysis, namely the de novo discovery of genes, the assembly of gene isoforms, and the accurate estimation of transcript abundance at the gene or the isoform level. Current RNA-Seq experiments are based on second-generation sequencing (SGS) instruments capable of generating a large number of short reads. From these reads, one obtains two types of information: (i) frequency of reads aligned to a contiguous genomic segment (exonic reads) and (ii) frequency of reads aligned to two contiguous segments of the genome with a single gap of from 60 bp to 400 kbp in size (junction reads) (9-11). If the set of possible isoforms is assumed known (i.e., the gene is well annotated), then it is possible to infer isoform-specific expression from exonic reads and junction reads based on simple statistical models such as the Poisson deconvolution model of Jiang and Wong (12). On the other hand, if the set of isoforms is not known or only partially known, then currently there is great

difficulty in isoform quantification based on SGS data. The main reason is due to insufficient length of the SGS reads. The median length of human gene transcripts is about 2,500 bp, which is much longer than the length of a contiguous read (about 250 bp) currently attainable by SGS. In previous work we showed that generally, isoform deconvolution from short-read RNA-Seq data is not an identifiable problem (13), in the sense that isoform expression cannot always be uniquely determined from the set of exons and splice junctions. Thus, strong assumptions are made on the set of candidate isoforms in all current methods for isoform assembly from short reads, including Cufflinks (14), SLIDE (15), and Montebello (16); as a result, the assembled isoforms are of uncertain accuracy. Although hundreds of RNA-Seq datasets are being generated in any given day by diverse groups in academia and industry, their interpretations all depend critically on the completeness and reliability of gene and isoform annotations on the species and cell types being analyzed. Although results from de novo transcript reconstruction algorithms can provide useful hints, they are not accurate enough to stand on their own as definitive evidence for new transcripts.

One may hope that the completeness and reliability of gene annotation are improving commensurably with the exponential increase in the amount and diversity of sequence data from RNA-Seq. However, from release 43 (September 2010) to 49

#### Significance

Isoform identification and discovery are an important goal for transcriptome analysis because the majority of human genes express multiple isoforms with context- and tissue-specific functions. Better annotation of isoforms will also benefit downstream analysis such as expression quantification. Current RNA-Seq methods based on short-read sequencing are not reliable for isoform discovery. In this study we developed a new method based on the combined analysis of short reads and long reads generated, respectively. <u>by reand and kind concertion</u> con-

quencing and applied characterization of t stem cell. The result specificity can be ac

Author contributions: KF.A. HV.B., R.A.R.-y. and J.G.U po new reagents/analytic tools; paper. The authors declare no confl Freely available online throu Data deposition: The data re pression Omnibus (GEO) data <sup>1</sup>To whom correspondence jundy@uw.edu.

cat 1073/pnas.1320101110/-/DCSupp



PACBIO

**8,048 RefSeq-annotated, full-length isoforms** and 5,459 predicted isoforms

"Over **one-third of these are novel isoforms**, including 273 RNAs from gene loci that have not previously been identified"

- <u>Au et al. (2013) Characterization of the human ESC transcriptome by hybrid sequencing. PNAS doi:</u> <u>10.1038/pnas.1320101110.</u>

#### NOVEL ISOFORMS IDENTIFIED IN RICE CULTIVARS WITH FULL-LENGTH ISO-SEQ™ SEQUENCING

ס- כן כל אכן כל איכן כל איכן כל איכן כל איכן כל איכ



# Transcript length distribution differences can be observed between sequencing platforms

PAG 2015 Poster: <u>Dario Copetti, et. al. "Rapid Full-Length Iso-Seq cDNA sequencing of Rice</u> mRNA to Facilitate Annotation and Identify Splice-Site Variation"

#### **ISO-SEQ<sup>™</sup> FOR IMPROVED ISOFORM DIFFERENTIATION**

לא ליק כלי כן כל יכן כל יכן כל יכן



Full-length transcript sequencing detects novel splice-site variants:

- alternative promoter/ poly(A)
- retained introns
- skipped exons
- alternative splice sites

PACBIO®

#### **DISCOVERING FUNGAL TRANSCRIPTOME DIVERSITY**



**Citation:** Gordon SP, Tseng E, Salamov A, Zhang J, Meng X, Zhao Z, et al. (2015) Widespread Fig 2. Long, high-quality, consensus sequences accurately benchmark transcript diversity.



Gordon SP, Tseng E, Salamov A, Zhang J, Meng X, et al. (2015) Widespread Polycistronic Transcripts in Fungi Revealed by Single-Molecule mRNA Sequencing. PLoS ONE 10(7): e0132628. doi:10.1371/journal.pone.0132628 http://journals.plos.org/plosone/article?id=info:doi/10.1371/journal.pone.0132628



PACBIO® **Epigenetics And Base modifications** 

### CHARACTERIZE EPIGENOMES UNIQUE TO SMRT® SEQUENCING



11450–11462 Nucleic Acids Research, 2012, Vol. 40, No. 22 doi:10.1093/nar/gks891 Published online 2 October 2012

#### The methylomes of six bacteria

lain A. Murray<sup>1</sup>, Tyson A. Clark<sup>2</sup>, Richard D. Morgan<sup>1</sup>, Matthew Boitano<sup>2</sup>, Brian P. Anton<sup>1</sup>, Khai Luong<sup>2</sup>, Alexey Fomenkov<sup>1</sup>, Stephen W. Turner<sup>2</sup>, Jonas Korlach<sup>2,\*</sup> and Richard J. Roberts<sup>1,\*</sup>

<sup>1</sup>New England Biolabs, 240 County Road, Ipswich, MA 01938 and <sup>2</sup>Pacific Biosciences, 1380 Willow Road, Menlo Park, CA 94025, USA



יס ארק כל ארכן כל ארכן כל ארכן כל ארכן כל ארכ<mark>ן כל ארכן כל ארכ</mark>



http://dx.doi.org/10.1111/j.1574-6976.2005.00006.x

#### HAEMOPHILUS INFLUENZAE EPIGENETICS



## in non-typeable Haemophilus influenzae

John M. Atack<sup>1</sup>, Yogitha N. Srikhanta<sup>1,†</sup>, Kate L. Fox<sup>2</sup>, Joseph A. Jurcisek<sup>3</sup>, Kenneth L. Brockman<sup>3</sup>, Tyson A. Clark<sup>4</sup>, Matthew Boitano<sup>4</sup>, Peter M. Power<sup>1</sup>, Freda E.-C. Jen<sup>1</sup>, Alastair G. McEwan<sup>2</sup>, Sean M. Grimmond<sup>5,†</sup>, Arnold L. Smith<sup>6</sup>, Stephen J. Barenkamp<sup>7</sup>, Jonas Korlach<sup>4</sup>, Lauren O. Bakaletz<sup>3</sup> & Michael P. Jennings<sup>1</sup>

Non-typeable *Haemophilus influenzae* contains an N<sup>6</sup>-adenine DNA-methyltransferase (ModA) that is subject to phase-variable expression (random ON/OFF switching). Five *modA* 

#### **EPIGENETIC SWITCH REGULATES ADAPTIVE MECHANISMS**

Table 1 | Summary of SMRT sequencing and methylome analysis of representative strains containing the five *modA* alleles under study.

modAallele	NTHi strain	Clinical symptoms	Accession number	Methylation sequence	Systematic name	Number of sites in genome	Genome size (bp)	Predicted ORFs
modA2	723	OM	CP007472	5'-CCGA( <sup>m6</sup> A)-3'	M.Hin723I	2,270	1,887,620	1,868
modA4	C486	OM	CP007471	5'-CG( <sup>m6</sup> A)G-3'	M.HinC486I	6,203	1,846,507	1,783
modA5	477	OM	CP007470	5'-AC( <sup>m6</sup> A)GC-3'	M.Hin477I	2,548	1,846,259	1,813
modA9	1209	OM	JMQP01000000	5'-CCTG( <sup>m6</sup> A)-3'	M.Hin1209I	2,504	1,895,979	2,247
modA10	R2866	Blood	CP002277*	5'-CCT( <sup>m6</sup> A)C-3'	M.Hin2866I	1,244	1,932,238	1,905

\*Strain already annotated and submitted (October 2010) to the EMBL/GenBank/DDBJ databases. A full summary of SMRT sequencing/methylome analysis derived data is presented in Supplementary Fig. 1.



ס- כן כל אכן כל אכן כל אכן כל איכן כל איכ

#### DNA METHYLATION IN C. ELEGANS

### Article

#### DNA Methylation on N<sup>6</sup>-Adenine in C. elegans

Eric Lieberman Greer,<sup>1,2,5,\*</sup> Mario Andres Blanco,<sup>1,2,5</sup> Lei Gu,<sup>1,2</sup> Erdem Sendinc,<sup>1,2</sup> Jianzhao Liu,<sup>3</sup> David Aristizábal-Corrales,<sup>1,2</sup> Chih-Hung Hsu,<sup>1,2</sup> L. Aravind,<sup>4</sup> Chuan He,<sup>3</sup> and Yang Shi<sup>1,2,\*</sup> <sup>1</sup>Division of Newborn Medicine, Children's Hospital Boston, 300 Longwood Avenue, Boston, MA 02115, USA <sup>2</sup>Department of Cell Biology, Harvard Medical School, Boston, MA 02115, USA <sup>3</sup>Department of Chemistry and Institute for Biophysical Dynamics, Howard Hughes Medical Institute, The University of Chicago, Chicago, IL 60637, USA <sup>4</sup>National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD 208943, USA <sup>5</sup>Co-first author <sup>\*</sup>Correspondence: eric.greer@childrens.harvard.edu (E.L.G.), yshi@hms.harvard.edu (Y.S.) http://dx.doi.org/10.1016/j.cell.2015.04.005

#### SUMMARY

Cell

In mammalian cells, DNA methylation on the fifth position of cytosine (5mC) plays an important role as an epigenetic mark. However, DNA methylation was considered to be absent in *C. elegans* because of the lack of detectable 5mC, as well as homologs of the cytosine DNA methyltransferases. Here, using multiple approaches, we demonstrate the presence of adenine N<sup>6</sup>-methylation (6mA) in *C. elegans* DNA. We further demonstrate that this modification ininheritance in C. elegans involves mutation of the histone H3 lysine 4 dimethyl (H3K4me2) demethylase spr-5 (Katz et al., 2009), which is an ortholog of the mammalian LSD1/KDM1A (Shi et al., 2004). The spr-5 mutant worms initially do not exhibit phenotypes; however, after successive generations lacking this demethylase, they display a progressively increased infertility. This fertility decline is concomitant with a global increase in the activating histone mark H3K4me2 and decline in the repressive histone mark H3K9me3 (Greer et al., 2014; Katz et al., 2009; Kerr et al., 2014; Nottke et al., 2011). Despite the fact that early- and late-generation spr-5 mutant worms should

#### GENOME DISTRIBUTION OF N6-ADENINE IN C. ELEGANS

ס-רן כל - רן כל - רן כל - רן כל - ר



- Important discoveries:
  - Identification of N6mA in C. elegans
  - Examination N6mA distribution
  - Investigation into the role of N6mA methylase and demethylase in epigenetic inheritance

Greer et al. Cell 2015.



# Introducing the Sequel<sup>™</sup> System

ס- כן כל ארכן כ

### SEQUEL SYSTEM THE SCALABLE PLATFORM FOR SMRT<sup>®</sup> SEQUENCING

- Based on proven SMRT Technology
- Increased capacity with 1 Million ZMWs/SMRT Cell
- Scalable throughput
- Decreased footprint and weight



#### 

### CONCLUSIONS

- Value of SMRT Sequencing
  - Longest available read lengths >10kb
  - High consensus accuracy
  - Uniform, unbiased coverage



- Generate high quality references
- Reveal the complexity of the transcriptome
- Accelerate new discoveries in epigenetics











# akieu@pacb.com

For Research Use Only. Not for use in diagnostics procedures. © Copyright 2015 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.