



PACBIO®



High-Quality *De Novo* Insect Genome Assemblies using PacBio Sequencing

Jonas Korlach, Sarah Kingan

October 5, 2016

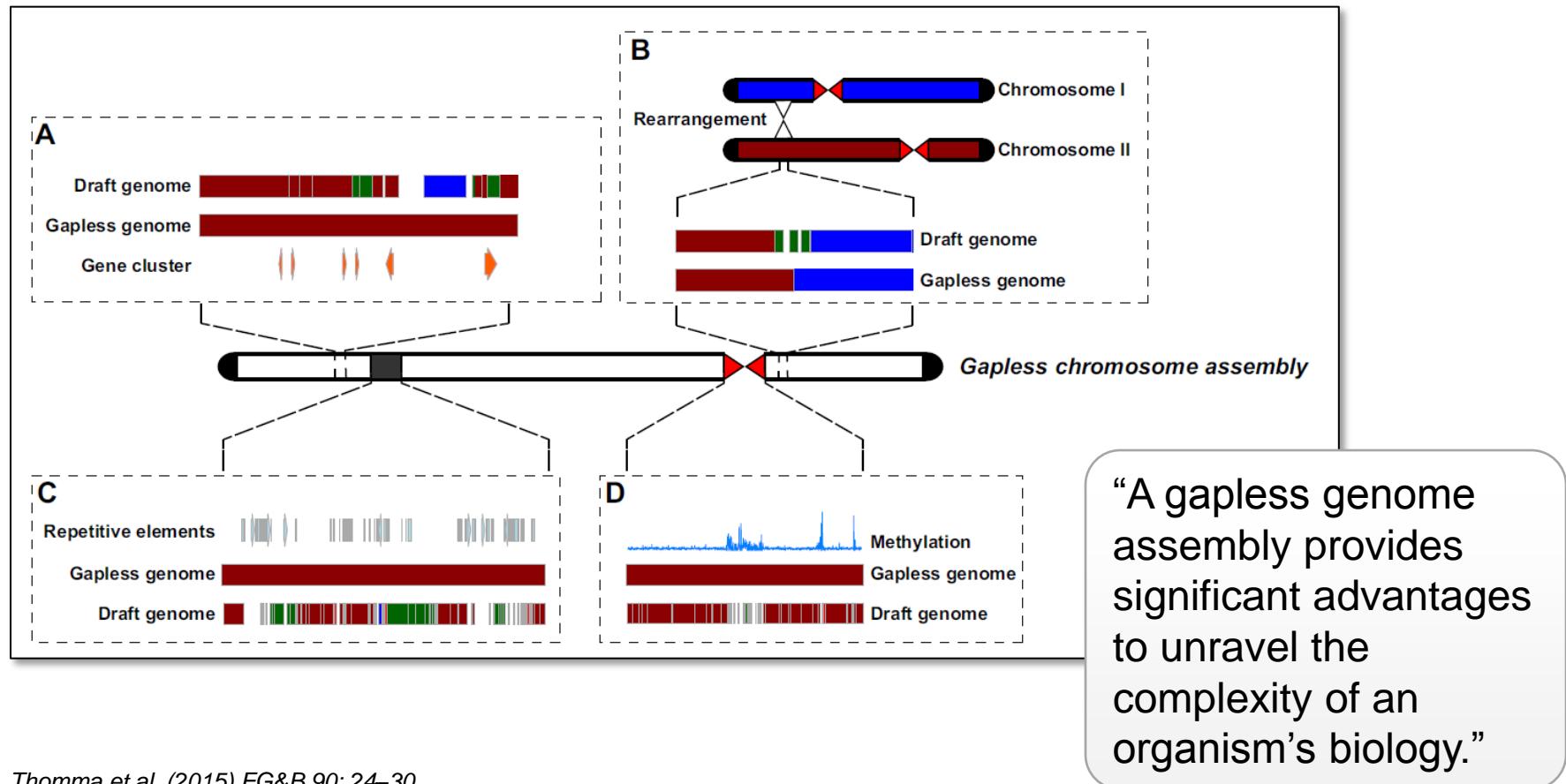
DRAFT vs. GAPLESS GENOMES

Mind the gap; seven reasons to close fragmented genome assemblies

Bart P.H.J. Thomma ^{a,*¹}, Michael F. Seidl ^a, Xiaoqian Shi-Kunne ^a, David E. Cook ^a, Melvin D. Bolton ^b, Jan A.L. van Kan ^a, Luigi Faino ^{a,¹}

^aLaboratory of Phytopathology, Wageningen University, Droevedaalsesteeg 1, 6708 PB Wageningen, The Netherlands

^bUnited States Department of Agriculture, Agricultural Research Service, Northern Crop Science Laboratory, Fargo, ND 58102-2765, USA



WHAT IS NEEDED FOR A HIGH-QUALITY GENOME?

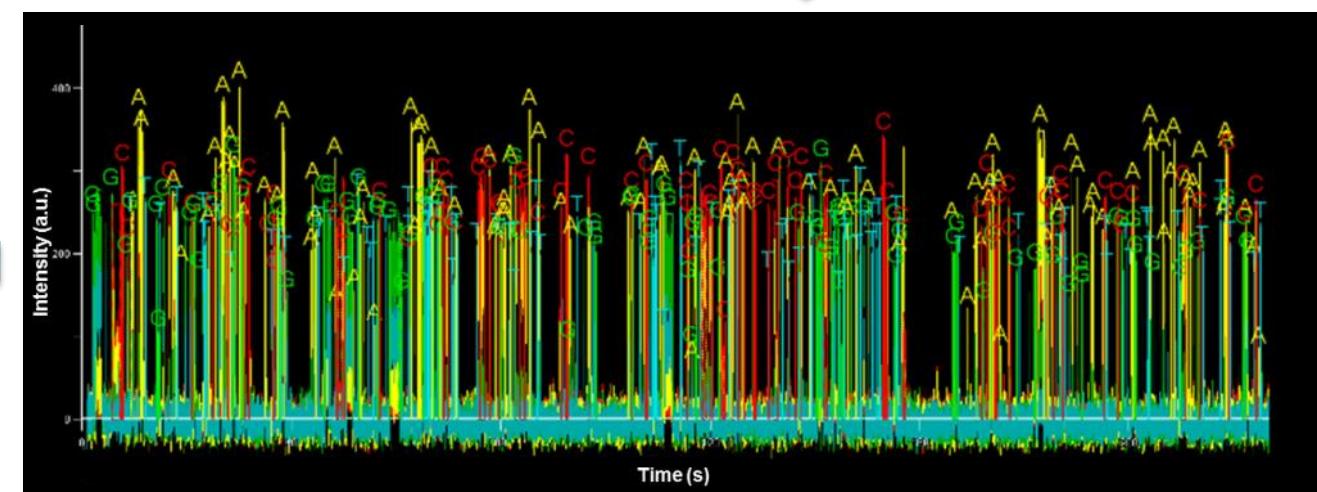
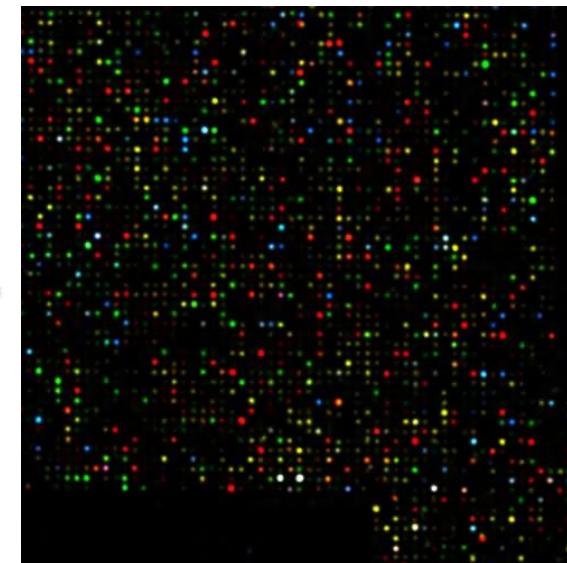
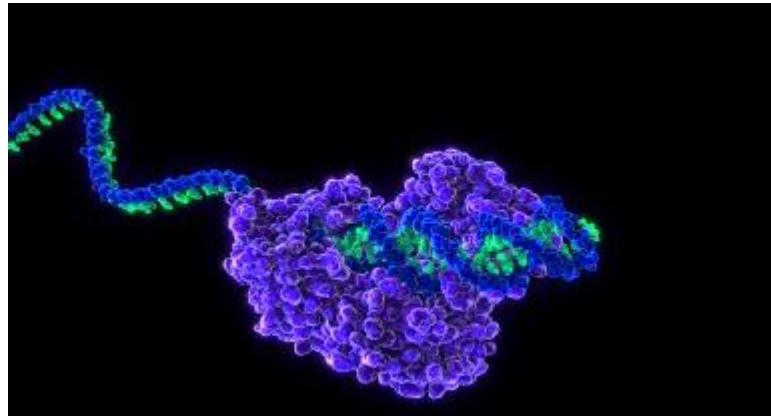
- Long sequence reads
 - Resolve repeats
 - Provides *contiguous* genome assembly
- Uniformity in sequencing coverage
 - Lack of bias (GC%, sequence complexity)
 - Allows sequencing the *entire* genome
- Absence of systematic sequencing errors
 - Random errors wash out in final consensus
 - Provides *accurate* genome sequence

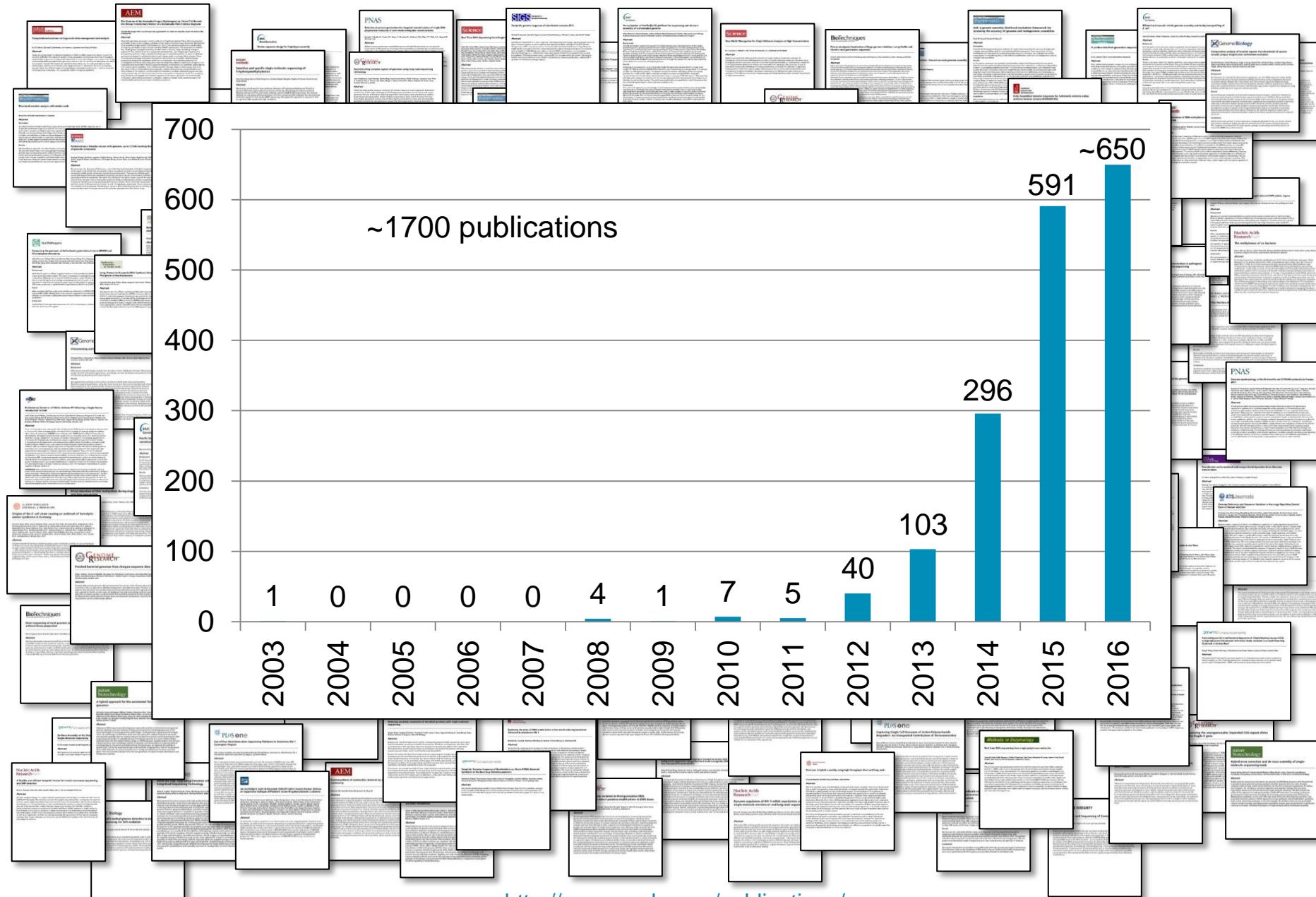
WHAT IS NEEDED FOR A HIGH-QUALITY GENOME?

- | | |
|---|--|
| <ul style="list-style-type: none">- Long sequence reads<ul style="list-style-type: none">- Resolve repeats- Provides <i>contiguous</i> genome assembly- Uniformity in sequencing coverage<ul style="list-style-type: none">- Lack of bias (GC%, sequence complexity)- Allows sequencing the <i>entire</i> genome- Absence of systematic sequencing errors<ul style="list-style-type: none">- Random errors wash out in final consensus- Provides <i>accurate</i> genome sequence | <p>PacBio¹</p> <p>average >10,000 bp</p> <p>least bias of any technology</p> <p>>99.999%</p> |
|---|--|

¹<http://www.pacb.com/wp-content/uploads/2015/09/Revolutionize-Genomics-with-SMRT-Sequencing.pdf>

SINGLE MOLECULE, REAL-TIME (SMRT) DNA SEQUENCING





HIGH-QUALITY PACBIO GENOME ASSEMBLIES



#1MbCtgClub

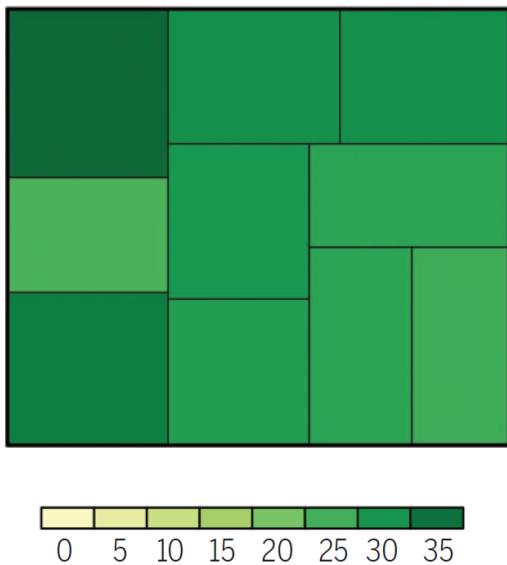
NEW GORILLA GENOME



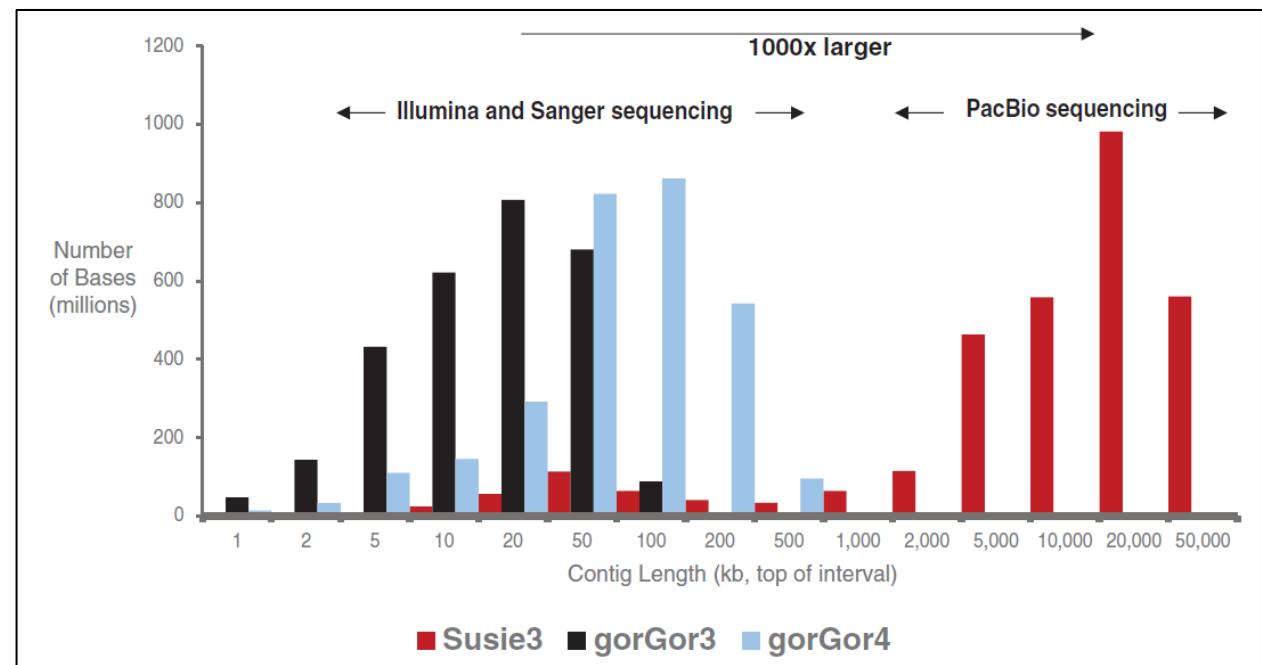
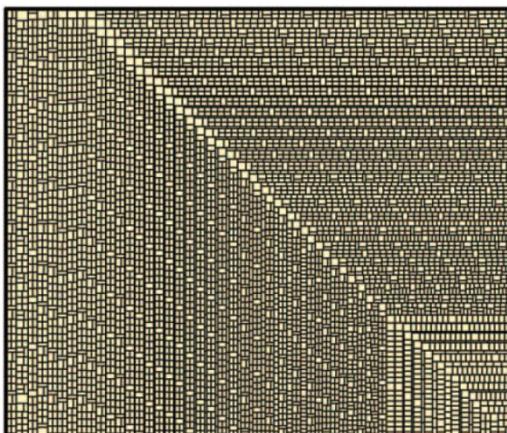
- >150fold greater contiguity
- Closure of >93% of previous gaps
- >148 Mb of new sequence
- Correction of previous errors
- 86% of structural variants are new
- Corrected previous estimates related to evolution

COMPARISON WITH PREVIOUS GORILLA ASSEMBLIES

B Long-read assembly (Susie3)



C Short-read assembly (gorGor3)



PACBIO PUBLICATIONS IN INSECT RESEARCH

- Bats may eat diurnal flies that rest on wind turbines
- Best practices in insect genome sequencing: What works and what doesn't.
- Birth of a new gene on the Y chromosome of *Drosophila melanogaster*.
- Buzz off, that's my bee!
- Comparative genome analysis of Wolbachia strain wAu
- Comparative genome sequencing reveals genomic signature of extreme desiccation tolerance in the anhydrobiotic midge.
- Complete genome of *Serratia* sp. strain FGI 94, a strain associated with leaf-cutter ant fungus gardens.
- Complete genome sequence of *Burkholderia* sp. strain RPE64, bacterial symbiont of the bean bug *Riptortus pedestris*.
- Complete genome sequence of *Endomicrobium proavitum*, a free-living relative of the intracellular symbionts of termite gut flagellates (phylum Elusimicrobia).
- Differing patterns of selection and geospatial genetic diversity within two leading *Plasmodium vivax* candidate vaccine antigens.
- DNA methylation on N6-adenine in *C. elegans*.
- Evolution of mosquito preference for humans linked to an odorant receptor.
- Genome analysis of a major urban malaria vector mosquito, *Anopheles stephensi*.
- Genome expansion via lineage splitting and genome reduction in the cicada endosymbiont *Hodgkinia*.
- Genome of *Cnaphalocrocis medinalis* granulovirus, the first Crambidae-infecting betabaculovirus isolated from rice leaffolder to sequenced.
- Genome sequence of the *Drosophila melanogaster* male-killing Spiroplasma strain MSRO endosymbiont.
- Genomes of 'Candidatus Liberibacter solanacearum' Haplotype A from New Zealand and the United States Suggest Significant Genome Plasticity in the Species.
- Genomics and host specialization of honey bee and bumble bee gut symbionts.
- Gut symbionts from distinct hosts exhibit genotoxic activity via divergent colibactin biosynthetic pathways.
- Insights into the preservation of the homomorphic sex-determining chromosome of *Aedes aegypti* from the discovery of a male-biased gene tightly linked to the M-locus.
- Landscape of standing variation for tandem duplications in *Drosophila yakuba* and *Drosophila simulans*.
- Large-scale mitogenomics enables insights into Schizophora (Diptera) radiation and population diversity.
- Long-read single molecule sequencing to resolve tandem gene copies: The Mst77Y region on the *Drosophila melanogaster* Y chromosome.
- Major improvements to the *Heliconius melpomene* genome assembly used to confirm 10 chromosome fusion events in 6 million years of butterfly evolution.
- *Microplitis demolitor* bracovirus proviral loci and clustered replication genes exhibit distinct DNA amplification patterns during replication.
- PacBio full-length transcriptome profiling of insect mitochondrial gene expression.
- *Paenibacillus* larvae-directed bacteriophage HB10c2 and its application in American Foulbrood-affected honey bee larvae.
- Quantitative profiling of *Drosophila melanogaster* Dscam1 isoforms reveals no changes in splicing after bacterial exposure.
- Selections that isolate recombinant mitochondrial genomes in animals.
- Single molecule long read sequencing resolves the detailed structure of complex satellite DNA loci in *Drosophila melanogaster*
- Site-specific genetic engineering of the *Anopheles gambiae* Y chromosome.
- Structural changes following the reversal of a Y chromosome to an autosome in *Drosophila pseudoobscura*
- Structure of the type IV secretion system in different strains of *Anaplasma phagocytophilum*.
- The bacterial microbiome of *Dermacentor andersoni* ticks influences pathogen susceptibility.
- The genome and methylome of a beetle with complex social behavior, *Nicrophorus vespilloides* (Coleoptera: Silphidae).
- The Glanville fritillary genome retains an ancient karyotype and reveals selective chromosomal fusions in Lepidoptera.
- The industrial melanism mutation in British peppered moths is a transposable element.
- The mitochondrial genome of a Texas outbreak strain of the cattle tick, *Rhipicephalus (Boophilus) microplus*, derived from whole genome sequencing Pacific Biosciences and Illumina reads.
- Unique transposon landscapes are pervasive across *Drosophila melanogaster* genomes.
- Variation and evolution in the glutamine-rich repeat region of *Drosophila argonaute-2*.
- Whole-genome sequence of *Bacillus* sp. SDL1, isolated from the social bee *Scaptotrigona depilis*.
- Whole-genome sequence of *Burkholderia* sp. strain RPE67, a bacterial gut symbiont of the bean bug *Riptortus pedestris*.
- Whole-genome sequence of *Erysipelothrix* larvae LV19(T) (=KCTC 33523(T)), a useful strain for arsenic detoxification, from the larval gut of the rhinoceros beetle, *Trypoxylus dichotomus*.
- Whole-genome sequence of *Serratia symbiotica* strain CWBI-2.3T, a free-living symbiont of the black bean aphid *Aphis fabae*.

DROSOPHILA SPECIES ASSEMBLIES

ARTICLES

**nature
biotechnology**

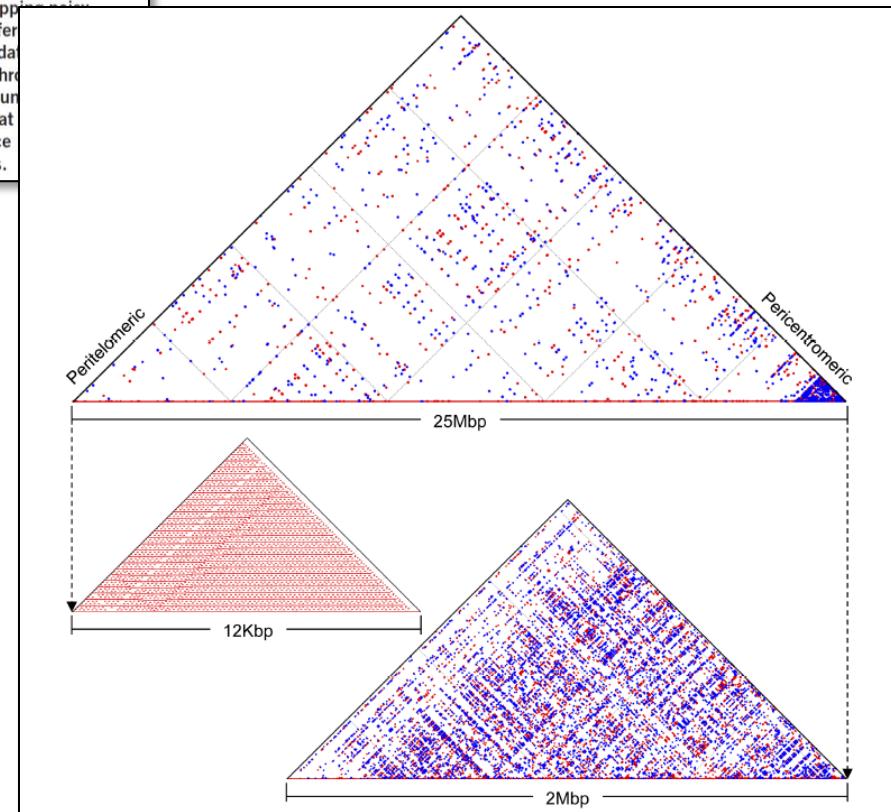
Assembling large genomes with single-molecule sequencing and locality-sensitive hashing

Konstantin Berlin^{1-3,6}, Sergey Koren^{4,6}, Chen-Shan Chin⁵, James P Drake⁵, Jane M Landolin⁵ & Adam M Phillippy⁴

Long-read, single-molecule real-time (SMRT) sequencing is routinely used to finish microbial genomes, but available assembly methods have not scaled well to larger genomes. We introduce the MinHash Alignment Process (MHAP) for overlapping long reads using probabilistic, locality-sensitive hashing. Integrating MHAP with the Celera Assembler enabled *de novo* assemblies of *Saccharomyces cerevisiae*, *Arabidopsis thaliana*, *Drosophila melanogaster* and a human hydatid cell line (CHM1) from SMRT sequencing. The resulting assemblies are highly continuous, include fully resolved chromosome arms and close persistent gaps in these reference genomes. Our assembly of *D. melanogaster* revealed previously unknown heterochromatic and telomeric transition sequences, and we assembled low-complexity sequences from CHM1 that were missing in the human GRCh38 reference. Using MHAP and the Celera Assembler, single-molecule sequencing can produce near-complete eukaryotic assemblies that are 99.99% accurate when compared with available reference genomes.

- Highly contiguous, including complete chromosome arms
- Filled persistent gaps in version 5 reference
- Revealed previously unknown heterochromatic and telomeric transition sequences

Single-contig assembly of *D. melanogaster* chromosome arm 3L



NEW INSIGHTS FROM PACBIO *DROSOPHILA* ASSEMBLIES

Long-Read Single Molecule Sequencing to Resolve Tandem Gene Copies: The *Mst77Y* Region on the *Drosophila melanogaster* Y Chromosome

Flavia J. Krsticevic,* Carlos G. Schrago,[†] and A. Bernardo Carvalho^{†,1}

^{*}Centro Internacional Franco Argentino de Ciencias de la Información y de Sistemas, CONICET, Ocampo y Esmeralda, S2000EPZ Rosario, Argentina, and [†]Departamento de Genética, Universidade Federal do Rio de Janeiro, 21941-971, Rio de Janeiro, Brazil

Nucleic Acids Research, 2015 1
doi: 10.1093/nar/gkv1193

Unique transposon landscapes are pervasive across *Drosophila melanogaster* genomes

Reazur Rahman¹, Gung-wei Chirn¹, Abhay Kanodia¹, Yuliya A. Sytnikova¹, Björn Brembs², Casey M. Bergman³ and Nelson C. Lau^{1,*}

¹Department of Biology and Rosenstiel Basic Medical Science Research Center, Brandeis University, Waltham, MA 02454, USA, ²Institute of Zoology, Universität Regensburg, Regensburg, Germany and ³Faculty of Life Sciences, University of Manchester, Manchester M21 0RG, UK

Birth of a new gene on the Y chromosome of *Drosophila melanogaster*

Antonio Bernardo Carvalho^{a,1}, Beatriz Vicoso^{a,b}, Claudia A. M. Russo^a, Bonnielin Swenor^c, and Andrew G. Clark^{c,1}

^aDepartamento de Genética, Universidade Federal do Rio de Janeiro, Caixa Postal 68011, 21941-971, Rio de Janeiro, Brazil; ^bInstitute of Science and Technology Austria, A-3400 Klosterneuburg, Austria; and ^cDepartment of Molecular Biology and Genetics, Cornell University, Ithaca, NY 14853

Variation and Evolution in the Glutamine-Rich Repeat Region of *Drosophila* Argonaute-2

William H. Palmer¹ and Darren J. Obbard

Institute of Evolutionary Biology and Centre for Infection, Evolution and Immunity, University of Edinburgh, EH9 3FL UK

Structural changes following the reversal of a Y chromosome to an autosome in *Drosophila pseudoobscura*

Ching-Ho Chang, Amanda M Larracuente

doi: <http://dx.doi.org/10.1101/058412>

Single molecule long read sequencing resolves the detailed structure of complex satellite DNA loci in *Drosophila melanogaster*

Daniel Emerson Khost, Danna G Eickbush, Amanda M Larracuente

doi: <http://dx.doi.org/10.1101/054155>

OVERVIEW

I. PacBio workflow

- Sample prep
- Sequencing
- Genome assembly

II. *Aedes aegypti* project

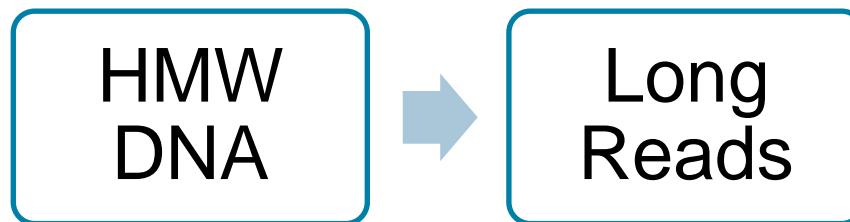
- Generating the PacBio assembly
- Quality and completeness
 - Assembly statistics
 - Phasing of allelic haplotypes
 - Full-length, base-pair resolution of insecticide-resistance gene

SIMPLE SCALABLE WORKFLOW

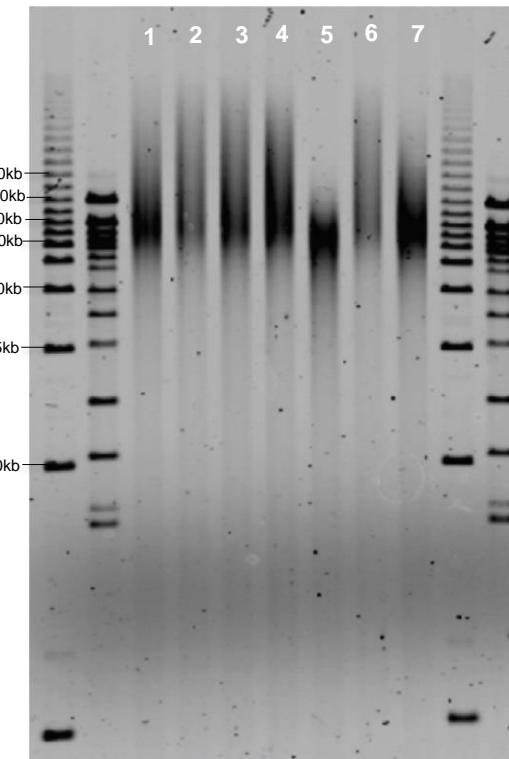
- SMRT Sequencing workflow is the same for all PacBio systems
 - Amplicons, RNA or large genomes



IMPORTANCE OF HIGH QUALITY TEMPLATE DNA



- DNA shearing
 - Megaruptor
 - Covaris g-TUBE
 - Needle shearing
 - 10 ug for 20 kb library
 - Multiple cells per library



DNA sheared with Megaruptor (Diagenode)

SMRTbell LIBRARY CONSTRUCTION



Template Preparation Workflow

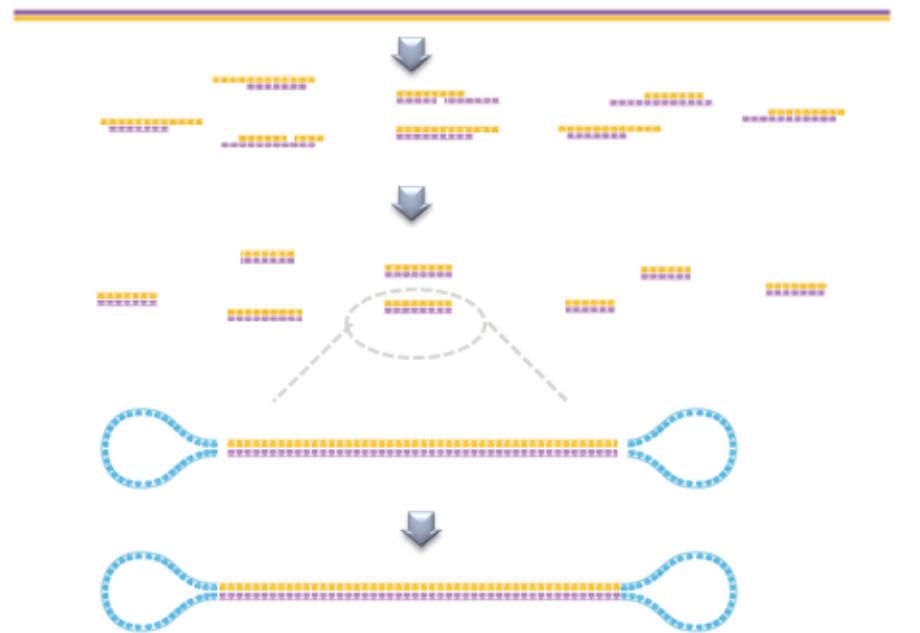
Fragment DNA and Concentration

DNA Damage Repair

Repair Ends

Ligate Adapters

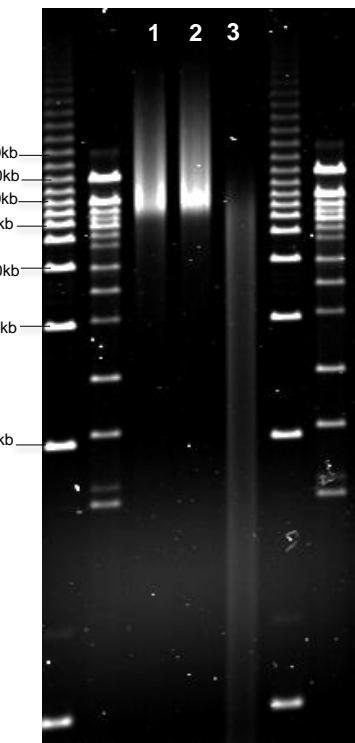
Purify Templates



Library construction takes approximately 3-4 hours

SIZE SELECTION OF SMRTbell LIBRARY

- Removes smaller library fragments
- Enables reads up to 60 kb long
- 20 – 30 kb size selection
- BluePippin or ELF (SageScience)



30 kb SMRTbell Library

ANNEAL PRIMER AND BIND POLYMERASE TO SMRTbell[®] TEMPLATES

LIBRARY CONSTRUCTION

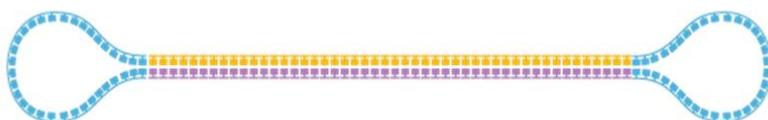
PRIMER ANNEALING AND POLY BINDING

LOADING

INSTRUMENT RUN

PRIMARY ANALYSIS

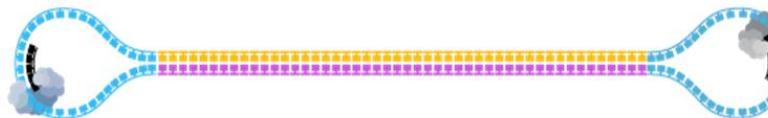
SECONDARY ANALYSIS



Anneal primer to both ends of the SMRTbell template



Bind polymerase to both ends of the SMRTbell template



PACBIO RS II AND SEQUEL SYSTEMS

	PacBio RS II	Sequel
Machine Launch	2013	2015
Number ZMW per Cell	150,000	1,000,000
Output per Cell	0.5 – 1 Gb	~ 5 Gb
Input DNA per Cell	~ 250 ng	~ 250 ng
N50 Read Length	15 – 20 kb	~15 kb
Consensus accuracy	QV50	QV50



LARGE GENOME SEQUENCING WITH SEQUEL

The screenshot shows a blog post titled "Sequel System Data Release: Arabidopsis Dataset and Genome Assembly" posted on Tuesday, September 27, 2016. The post discusses the release of the first *Arabidopsis thaliana* (Ler-0) dataset and genome assembly using the Sequel System. It highlights significant improvements in chemistry and reduced DNA requirements compared to previous releases. The PacBio logo is visible in the top left corner of the page.

PACBIO

PRODUCTS + SERVICES RESEARCH FOCUS APPLICATIONS SMRT SCIENCE SUPPORT COMPANY

BLOG

SMRT SEQUENCING

SMRT RESOURCES

- Scientific Publications
- Conference Proceedings
- PacBio Literature
- Video Gallery
- [Blog](#)

Sequel System Data Release: Arabidopsis Dataset and Genome Assembly

Tuesday, September 27, 2016

Today we are pleased to release the first *Arabidopsis thaliana* (Ler-0) dataset and *de novo* genome assembly generated with the Sequel System, using two SMRT Cells and 12 hours of runtime. Only three years ago, we released our first [genome assembly](#)¹ for *Arabidopsis* produced on the PacBio RS II using P4-C2 chemistry, 85 SMRT Cells and 255 hours of runtime. Four months later, we released a second [Arabidopsis dataset](#)¹ using the improved P5-C3 chemistry, which reduced the number of SMRT Cells to 46 and runtime to 138 hours.

We produced this Sequel dataset using our latest chemistry enhancements which significantly reduce the amount of DNA required. Prior to these chemistry improvements, the amount of DNA needed to run many large genome projects on the Sequel System was prohibitive. These modifications enable the use of loading concentrations equivalent to PacBio RS II levels.

LARGE GENOME SEQUENCING WITH SEQUEL

The screenshot shows a blog post on the PacBio website. The header features the PacBio logo and navigation links for Products + Services, Research Focus, Applications, SMRT Science, Support, and Company. A sidebar on the left lists "SMRT SEQUENCING" and "SMRT RESOURCES" (Scientific Publications, Conference Proceedings, PacBio Literature, Video Gallery, Blog). The main content area displays a photograph of a microscope. The blog post title is "Sequel System Data Release: Arabidopsis Dataset and Genome Assembly", dated Tuesday, September 27, 2016. The text describes the release of the first *Arabidopsis thaliana* (Ler-0) dataset and genome assembly using the Sequel System, generated with two SMRT Cells and 12 hours of runtime. It compares this to previous releases using PacBio RS II, noting significant improvements in both speed and DNA requirements.

Sequel System Data Release: Arabidopsis Dataset and Genome Assembly

Tuesday, September 27, 2016

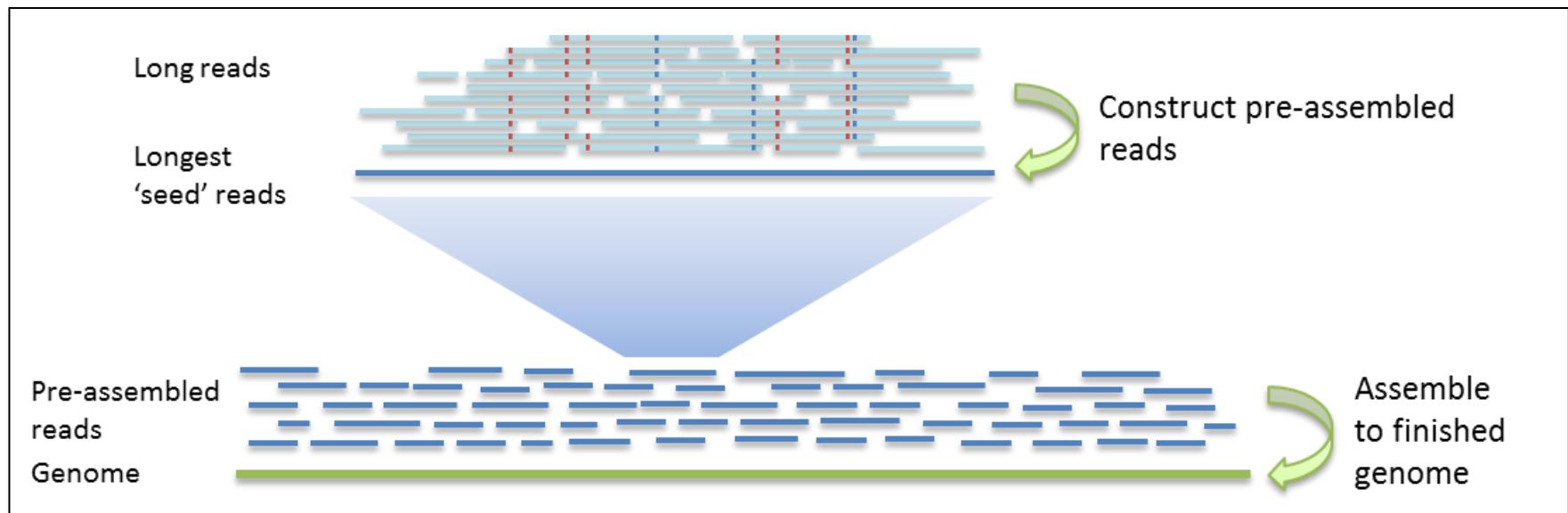
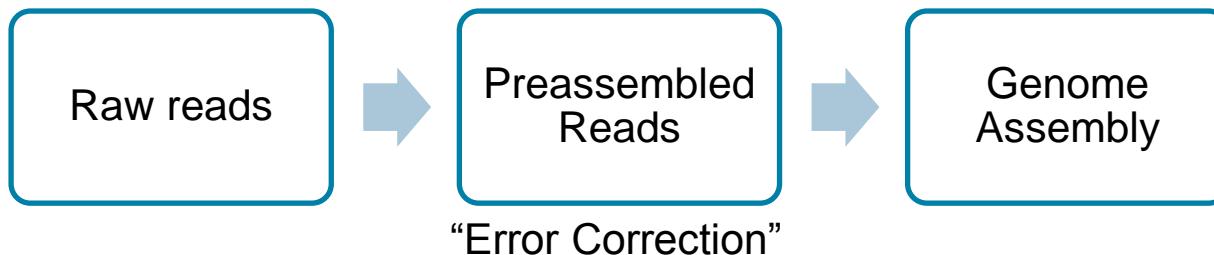
Today we are pleased to release the first *Arabidopsis thaliana* (Ler-0) dataset and *de novo* genome assembly generated with the Sequel System, using two SMRT Cells and 12 hours of runtime. Only three years ago, we released our first [genome assembly](#)¹ for *Arabidopsis* produced on the PacBio RS II using P4-C2 chemistry, 85 SMRT Cells and 255 hours of runtime. Four months later, we released a second [Arabidopsis dataset](#)¹ using the improved P5-C3 chemistry, which reduced the number of SMRT Cells to 46 and runtime to 138 hours.

We produced this Sequel dataset using our latest chemistry enhancements which significantly reduce the amount of DNA required. Prior to these chemistry improvements, the amount of DNA needed to run many large genome projects on the Sequel System was prohibitive. These modifications enable the use of loading concentrations equivalent to PacBio RS II levels.

Genome Length	120 Mb
Total BP	10.8 Gb (90-fold)
Number Cells	2
Read N50	16.4 kb
Input per Cell	0.25 ug

Assembly Length	122.9 Mb
Number Contigs	238
Contig N50	10.4 Mb
Max Contig	15.0 Mb

HIERARCHICAL GENOME ASSEMBLY PROCESS (HGAP)



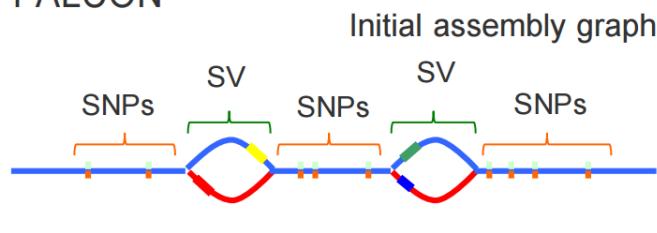
DIPLOID ASSEMBLY WITH FALCON-UNZIP



Jason Chin, PacBio

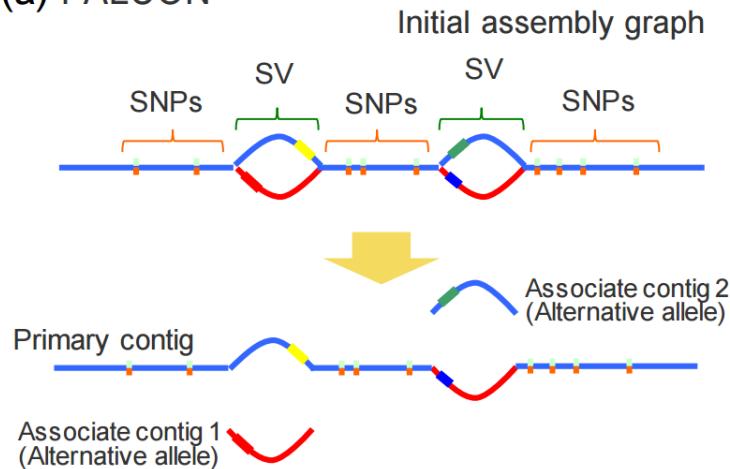
DIPLOID ASSEMBLY WITH FALCON-UNZIP

(a) FALCON



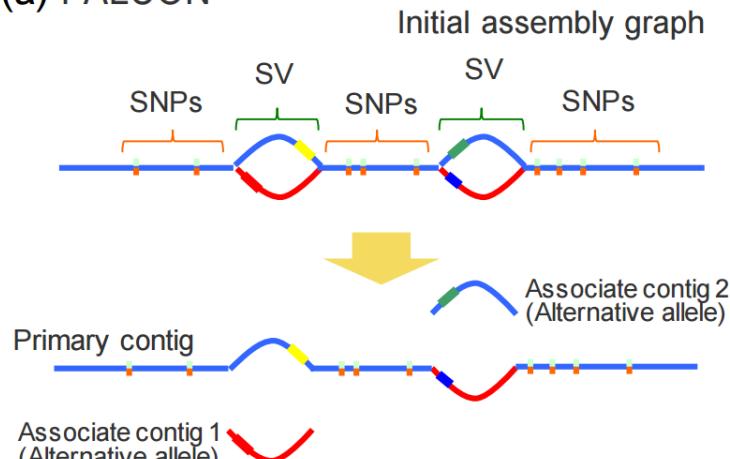
DIPLOID ASSEMBLY WITH FALCON-UNZIP

(a) FALCON



DIPLOID ASSEMBLY WITH FALCON-UNZIP

(a) FALCON



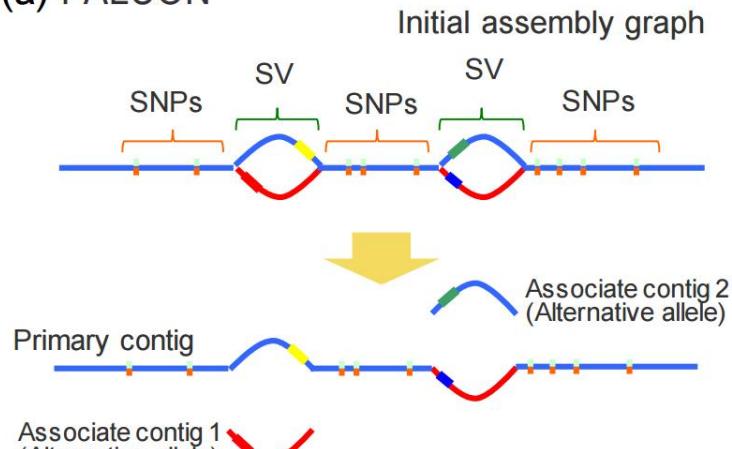
(b)



Phase heterozygous SNPs and identify the haplotype of each read

DIPLOID ASSEMBLY WITH FALCON-UNZIP

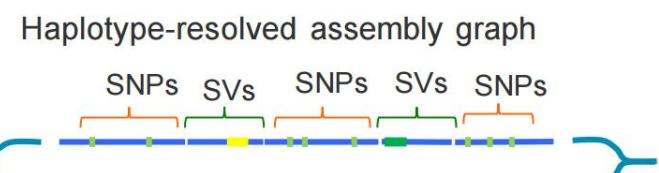
(a) FALCON



(b)

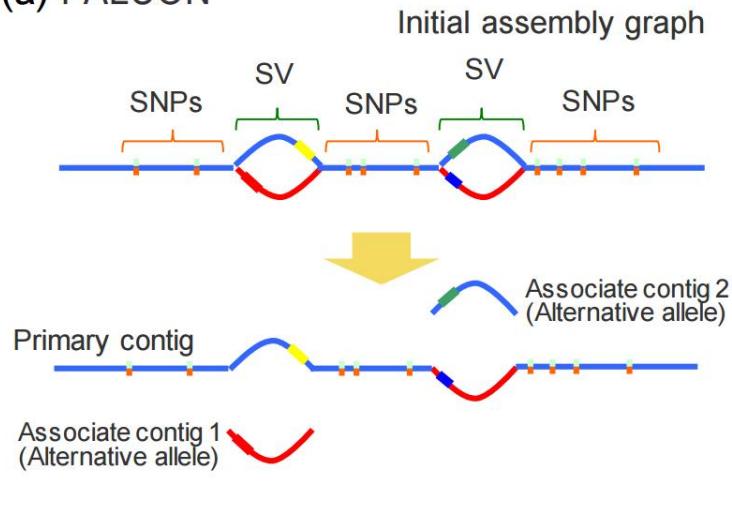


(c) FALCON-Unzip



DIPLOID ASSEMBLY WITH FALCON-UNZIP

(a) FALCON

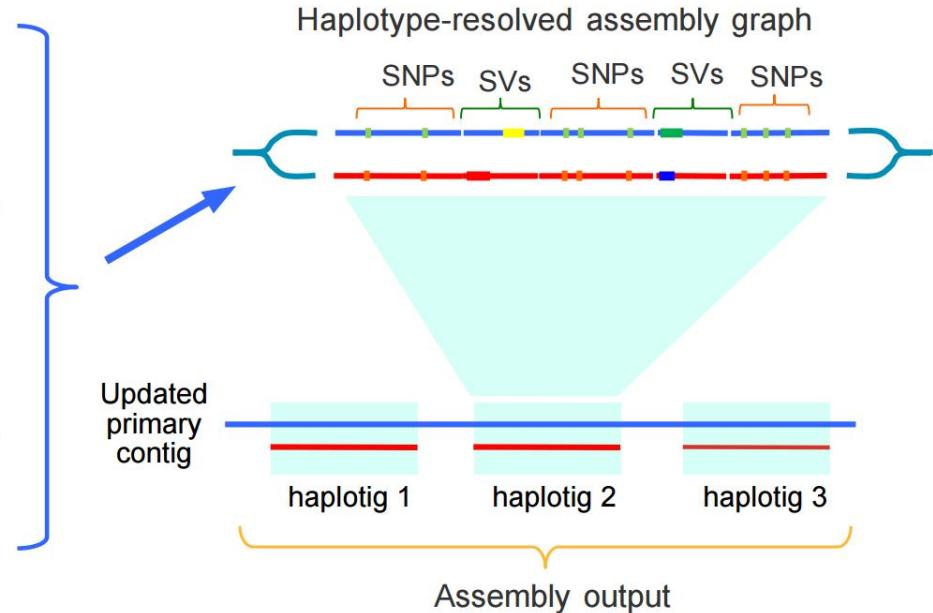


(b)



Phase heterozygous SNPs and
identify the haplotype of each read

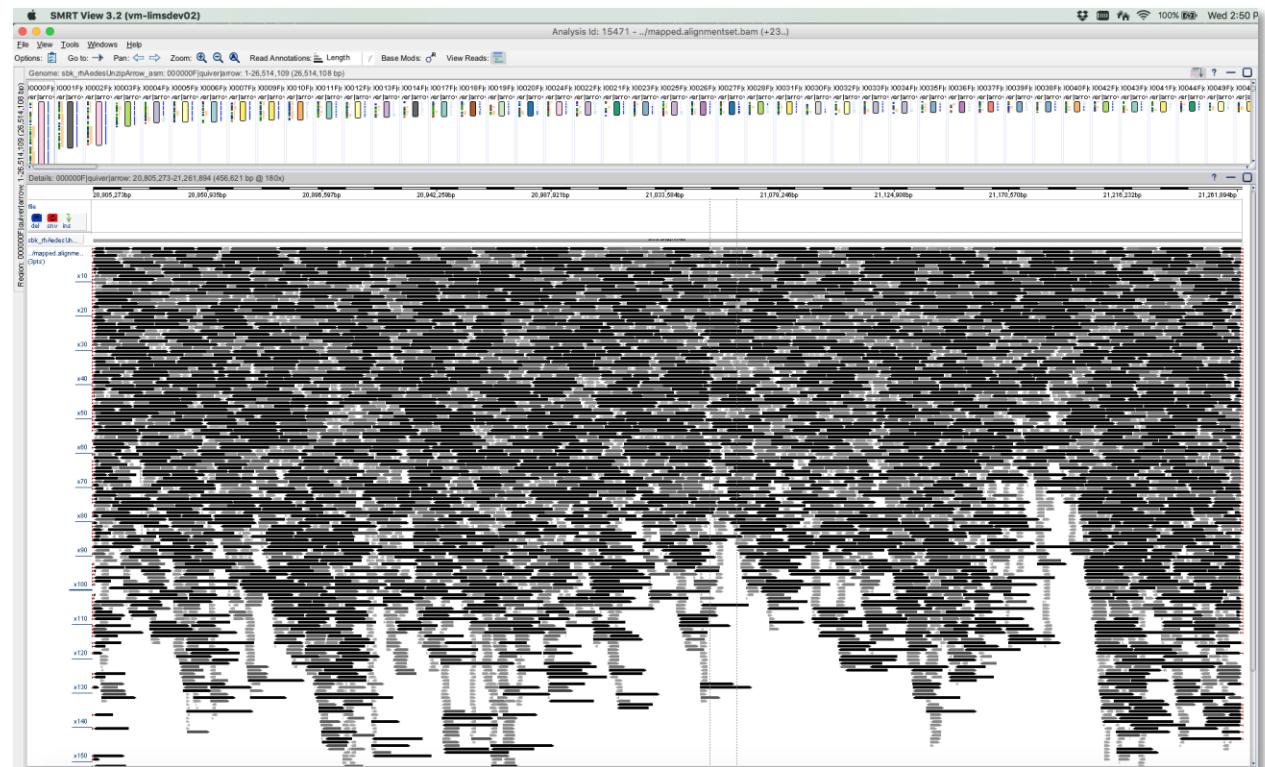
(c) FALCON-Unzip



GENOME POLISHING

Quiver and Arrow Algorithms

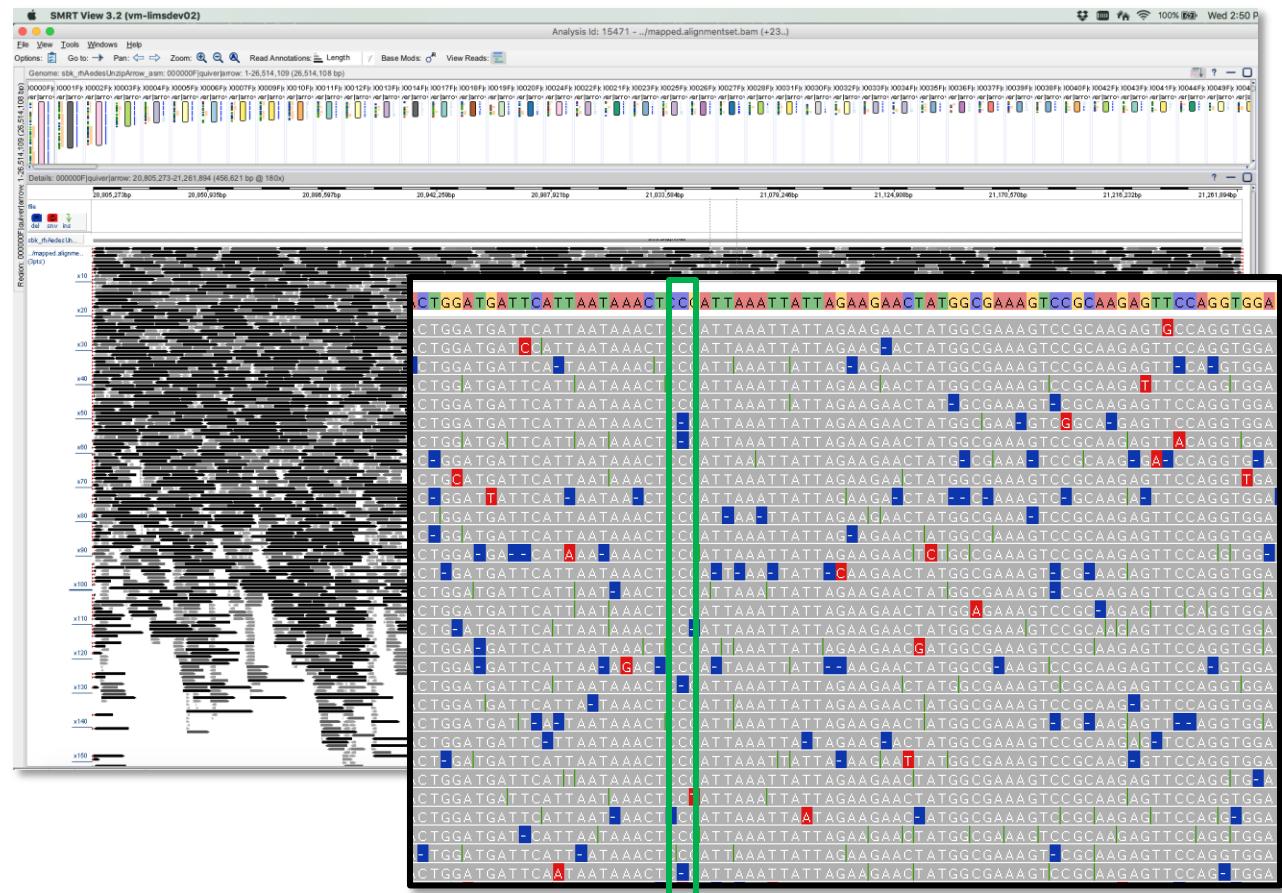
- Map raw reads back to genome sequence
- Compute consensus base and base quality
- Hidden Markov Model, trained on sequencing chemistry characteristics



GENOME POLISHING

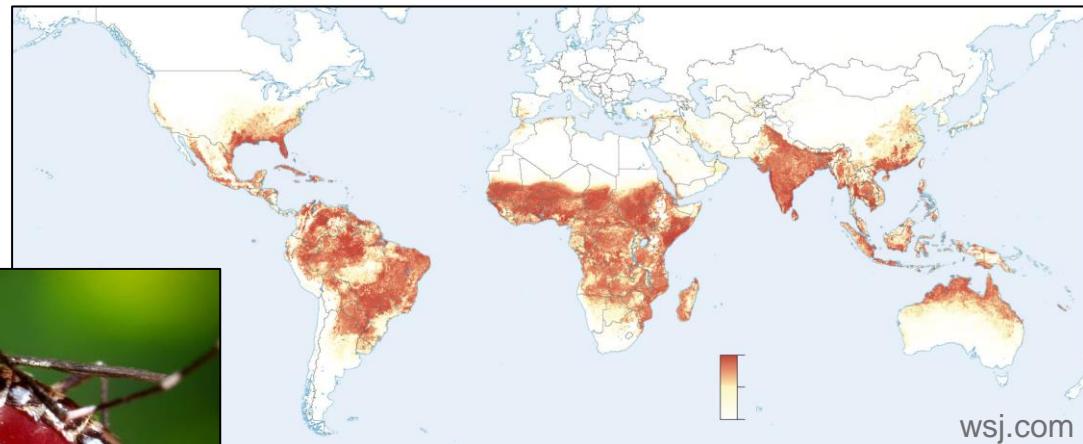
Quiver and Arrow Algorithms

- Map raw reads back to genome sequence
- Compute consensus base and base quality
- Hidden Markov Model, trained on sequencing chemistry characteristics



AEDES AEGYPTI GENOME WORKING GROUP (AGWG)

- Led by Drs. Leslie Vosshall and Ben Matthews, HHMI, Rockefeller University
- Dozens of academic researchers and private companies

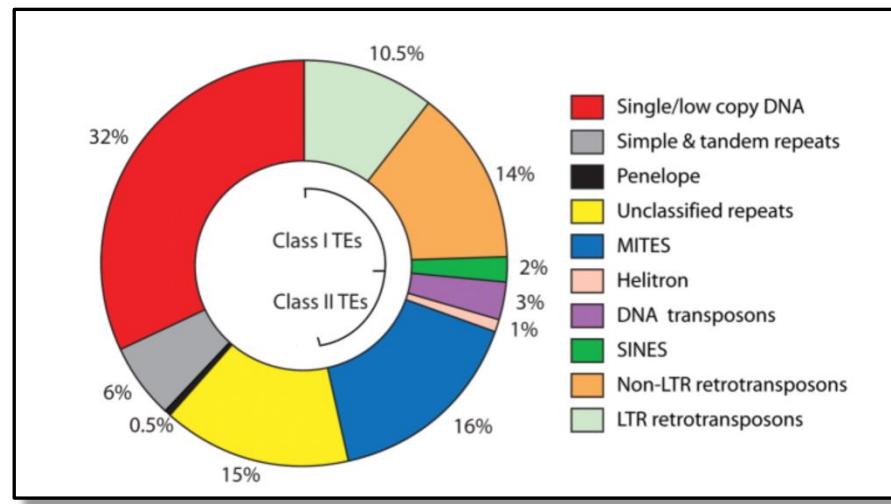


- Vector for yellow fever, Zika, dengue and other diseases
- Tropical and semi-tropical

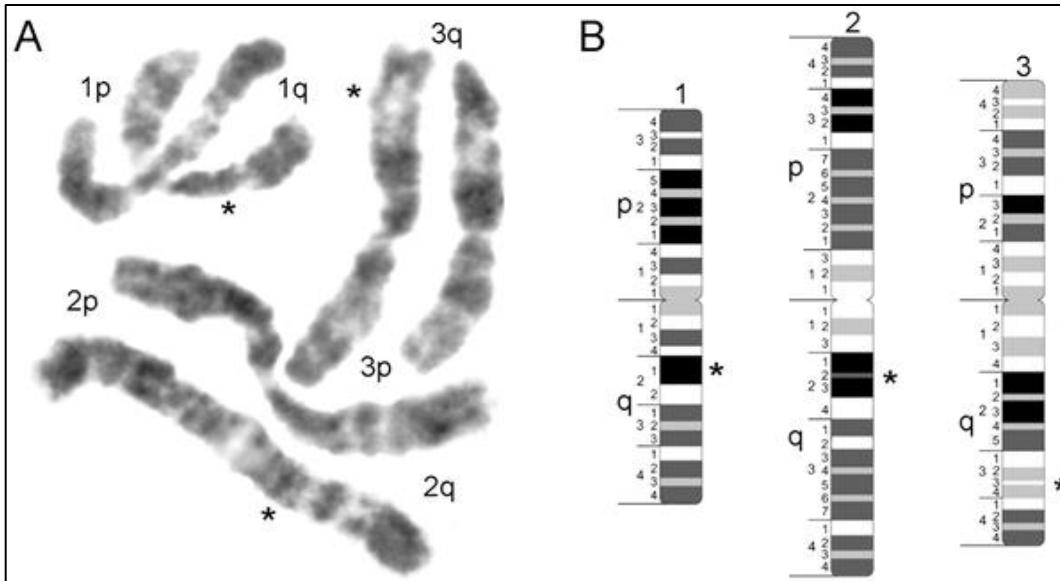
- Ancestral subspecies: *Aedes aegypti formosus* from Africa
- Anthropophilic *A. a. aegypti* dispersed by humans
- High selective pressure due to eradication efforts

AEDES AEGYPTI GENOME

- 3 metacentric chromosomes
- Homomorphic sex chromosomes
- 1.3 Gb
- 5 times the length of *Anopheles gambiae*
- Highly repetitive



Nene et al. 2007



Timoshevskiy et al. 2013

CURRENT REFERENCE



Nene et al. 2007

- Published in June, 2007
 - Sanger-based assembly of Liverpool inbred strain
 - Collaborative effort of private and public entities
 - Remains fragmented and low quality

AEDES AEGYPTI PACBIO ASSEMBLY

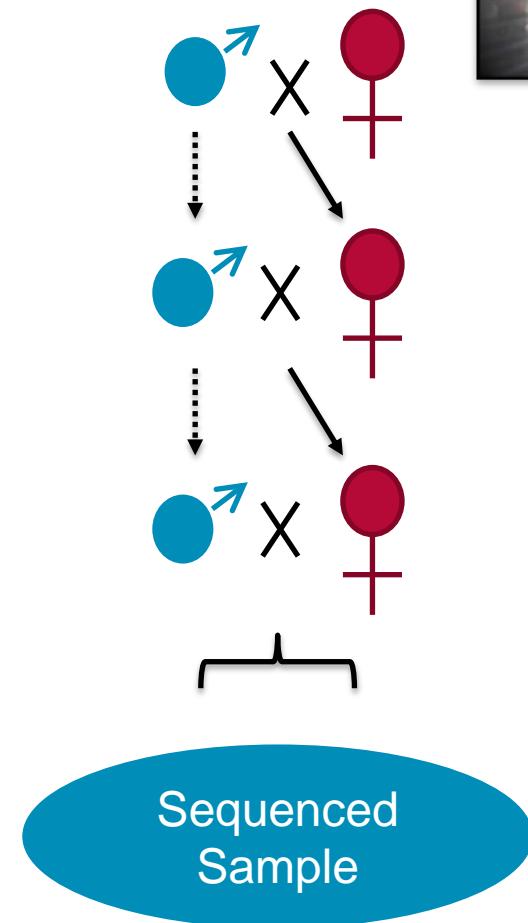
Overview

- Biological Sample
- Library Preparation
- Sequencing
- Diploid Genome Assembly
- Assessment of Genome Quality and Completeness
- Long-range Haplotype Resolution
- Full-length, Base-Pair Resolution of Insecticide Resistance Gene

BIOLOGICAL SAMPLE

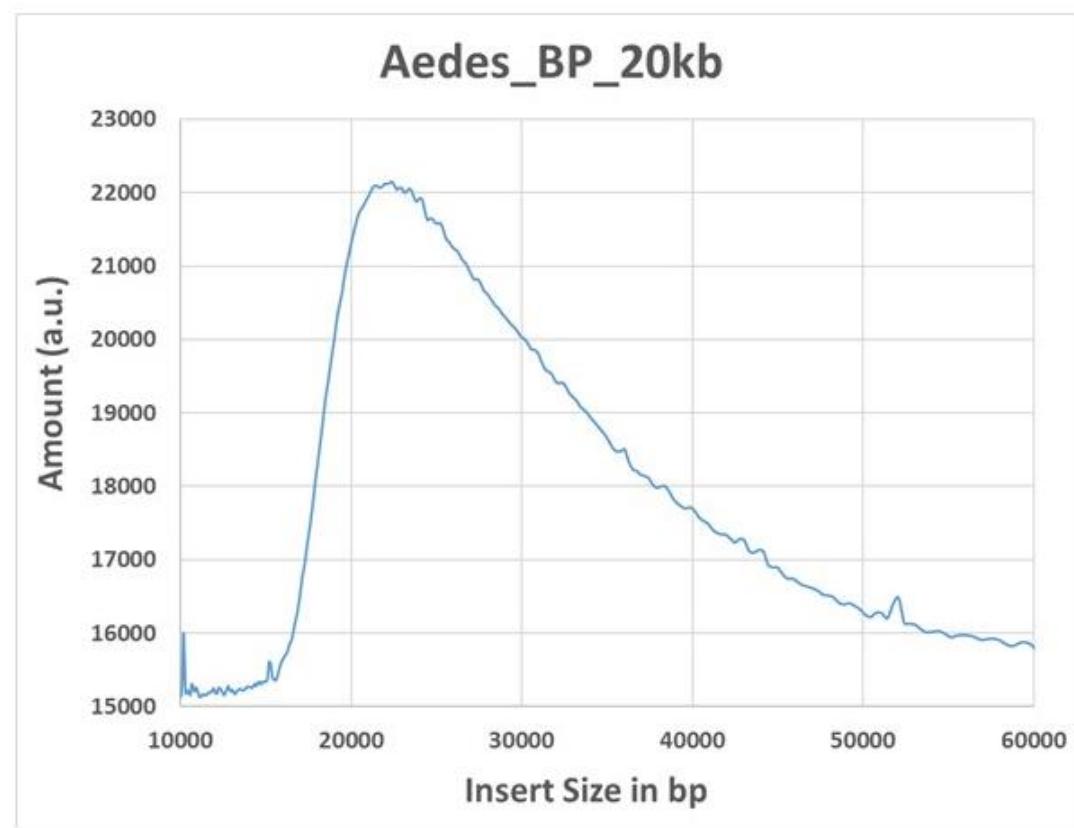
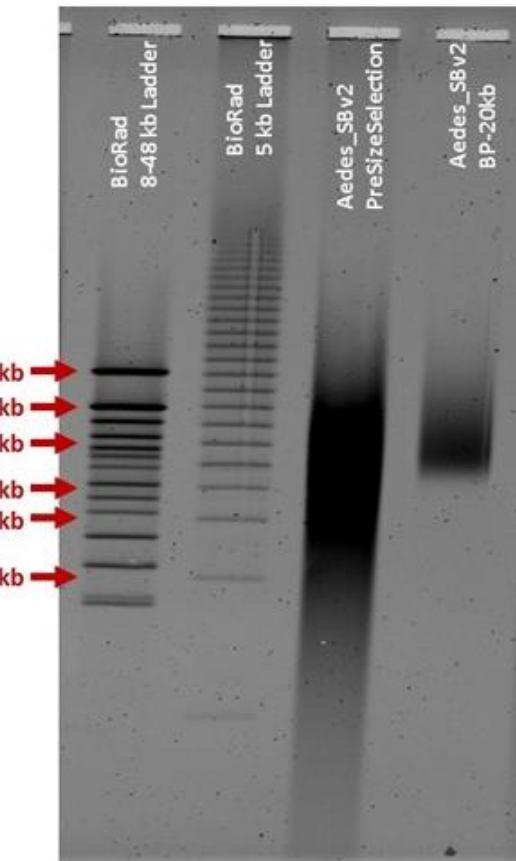
Ben Matthews, Rockefeller University

- Inbred (12 generation) strain matching current reference (“Liverpool”) is lost.
- New inbred strain generated:
 - Single male and female from founder strain
 - Male mated to single direct female descendant for 4 generations.
 - Ploidy $\leq 4N$
- 80 animals (male pupae)
- HMW genomic DNA extracted with MagAttract Kit, Qiagen



LONG-INSERT LIBRARIES

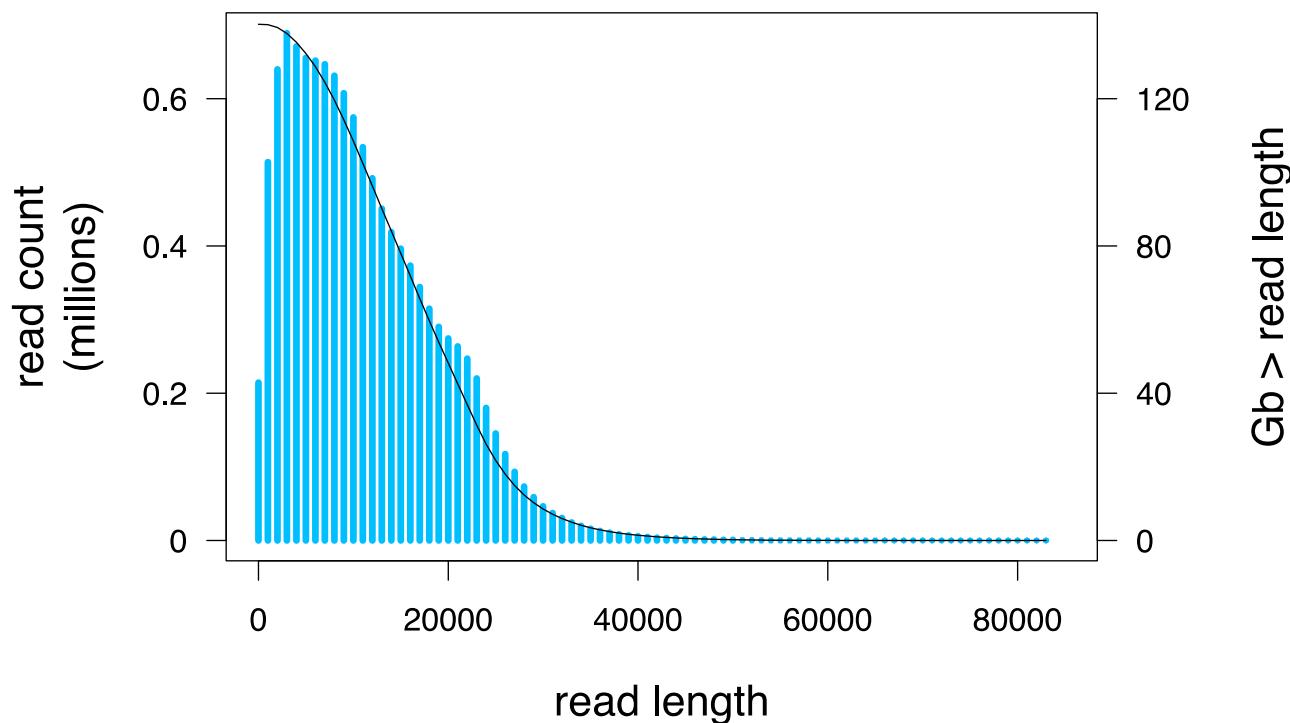
- Three libraries generated by Paul Peluso (PacBio)
- Size selection of SMRTbell library using BluePippin and ELF



TOTAL RAW READ DATA

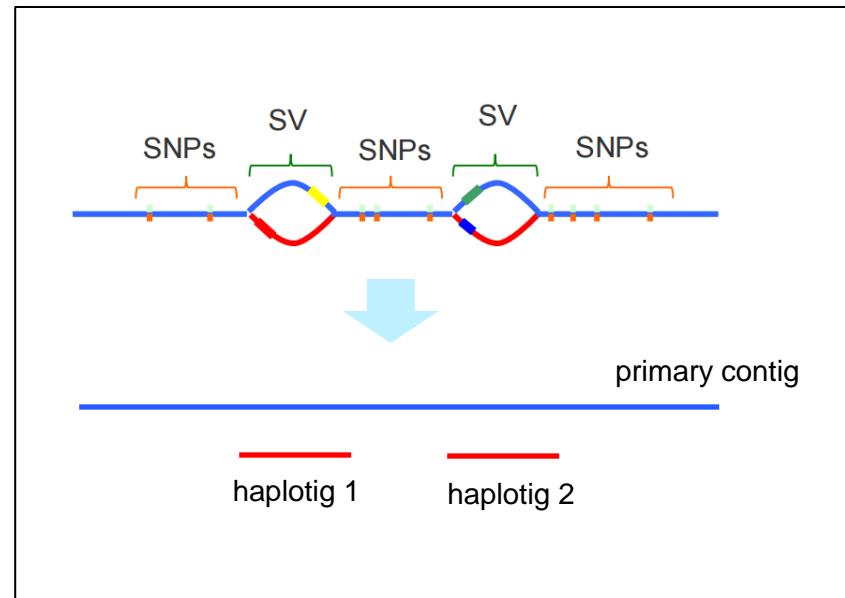
PacBio RS II

Total SMRT Cells	Total Bases	Genome Coverage	Subread N50
177	140 Gb	~108-fold	17 kb



DIPLOID ASSEMBLY

- Falcon-UNZIP assembly using p-reads longer than 19 kb
- Polishing with Arrow to QV40
- Identification of additional haplotigs using coverage and gene annotation



	Primary Contigs	Alternate Haplotigs
Total Length	1.45 Gb	594 Mb
Number of Contigs	3,462	4,328
Contig N50	1.43 Mb	382 kb
Longest Contig	26.5 Mb	5.42 Mb

PACBIO ASSEMBLY: IMPROVED QUALITY AND CONTINUITY

Assembly	L3.31 reference ¹	PacBio
Genome Size	1.38 Gb	1.45 Gb
Number of Contigs	36,204	3,462
Improvement factor		10
Contig N50	0.083 Mb	1.43 Mb
Improvement factor		17
CEGMA Genes	91%	92%
CEGMA Genes with frame shift errors	6 (2.7%)	4 (1.7%)

¹ftp://ftp.ensemblgenomes.org/pub/metazoa/release-31/fasta/aedes_aegypti/dna/

PACBIO ASSEMBLY: IMPROVED QUALITY AND CONTINUITY

Assembly	2015 Illumina ¹	PacBio
Genome Size	0.74 Gb	1.45 Gb
Number of Contigs	961,292	3,462
Improvement factor		278
Contig N50	0.001 Mb	1.43 Mb
Improvement factor		1430
CEGMA Genes	72%	92%
CEGMA Genes with frame shift errors	7 (3.9%)	4 (1.7%)

¹[ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA_001014885.1_ASM101488v1](http://ftp.ncbi.nlm.nih.gov/genomes/all/GCA_001014885.1_ASM101488v1)

LONG-RANGE HAPLOTYPE RESOLUTION

Primary Contig

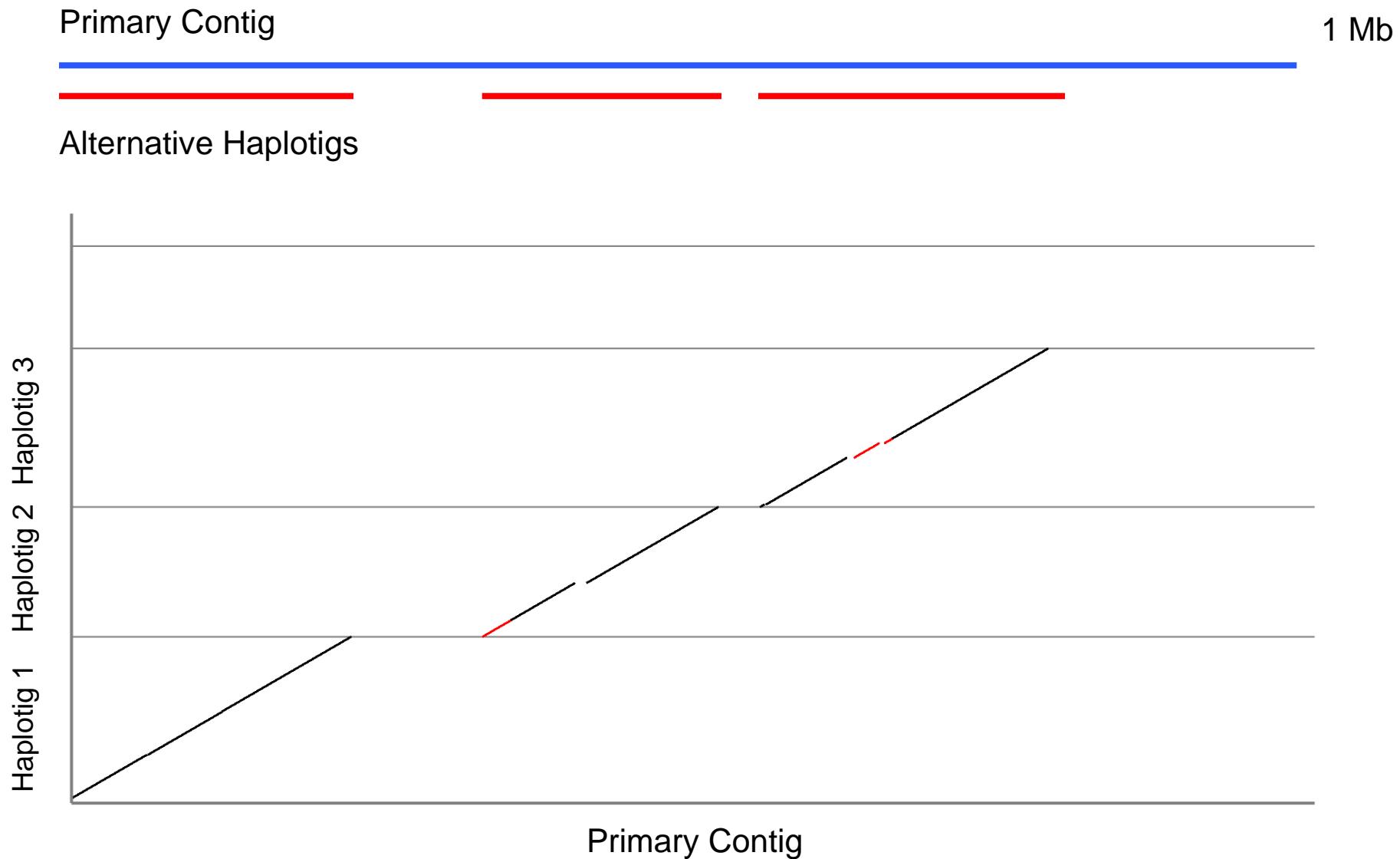


Alternative Haplotigs



1 Mb

LONG-RANGE HAPLOTYPE RESOLUTION



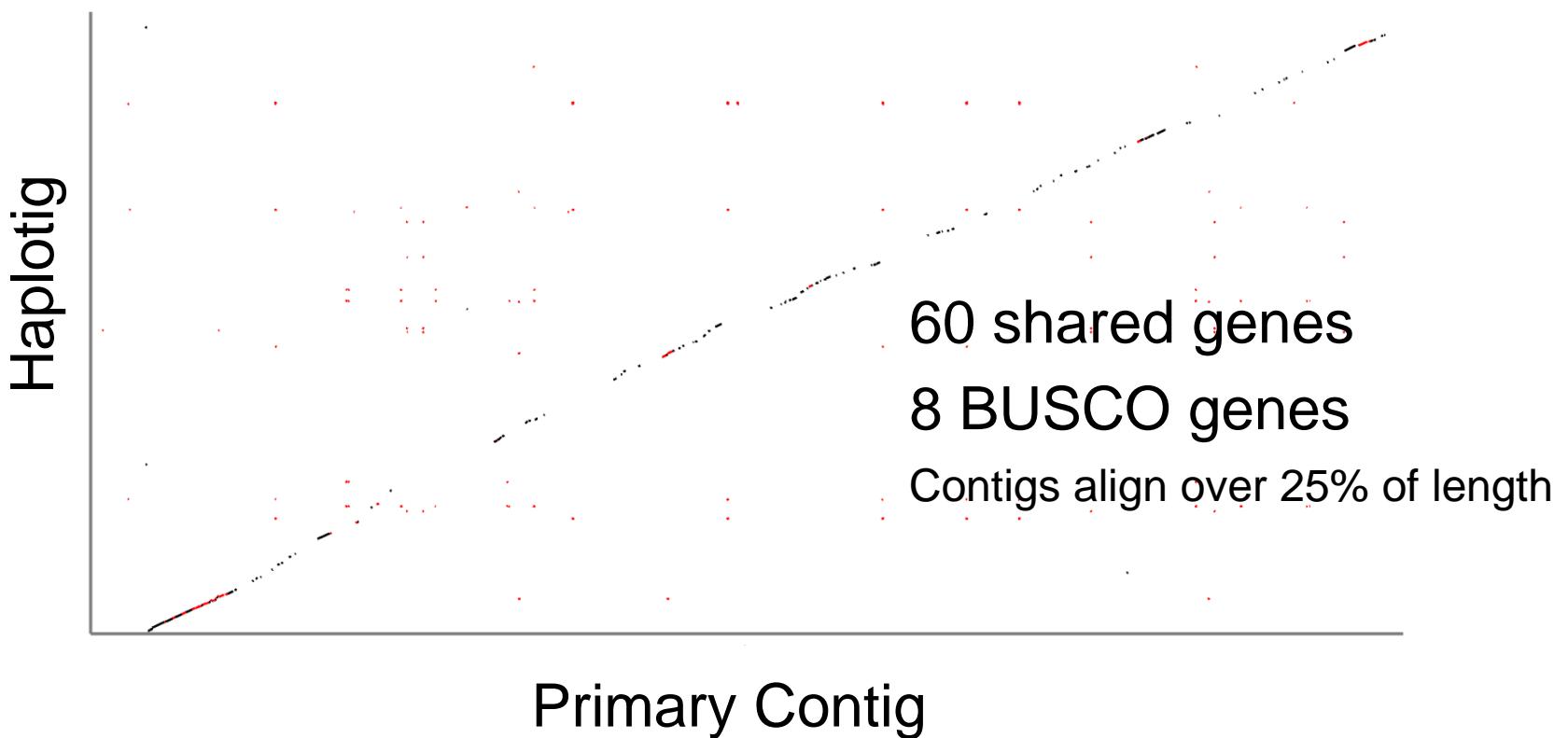
LONG-RANGE HAPLOTYPE RESOLUTION

Primary Contig



Alternative Haplotype

4 Mb



EXAMPLE - VOLTAGE-GATED SODIUM CHANNEL

- Target for insecticides, mutations likely involved in insecticide resistance¹

	PacBio assembly	L3.31 Reference ²	2015 Illumina assembly ³
Location	Contig 88 (2.87 Mb)	Scaffold 186 (1.96 Mb)	Scaffold 76061 (0.011 Mb)
Gene Status	Complete ~500 kb gene size 2109 a.a. protein	Incomplete 11% of protein is missing (N-terminal 225 a.a.)	Incomplete 84% of protein is missing (N-terminal 1766 a.a.)

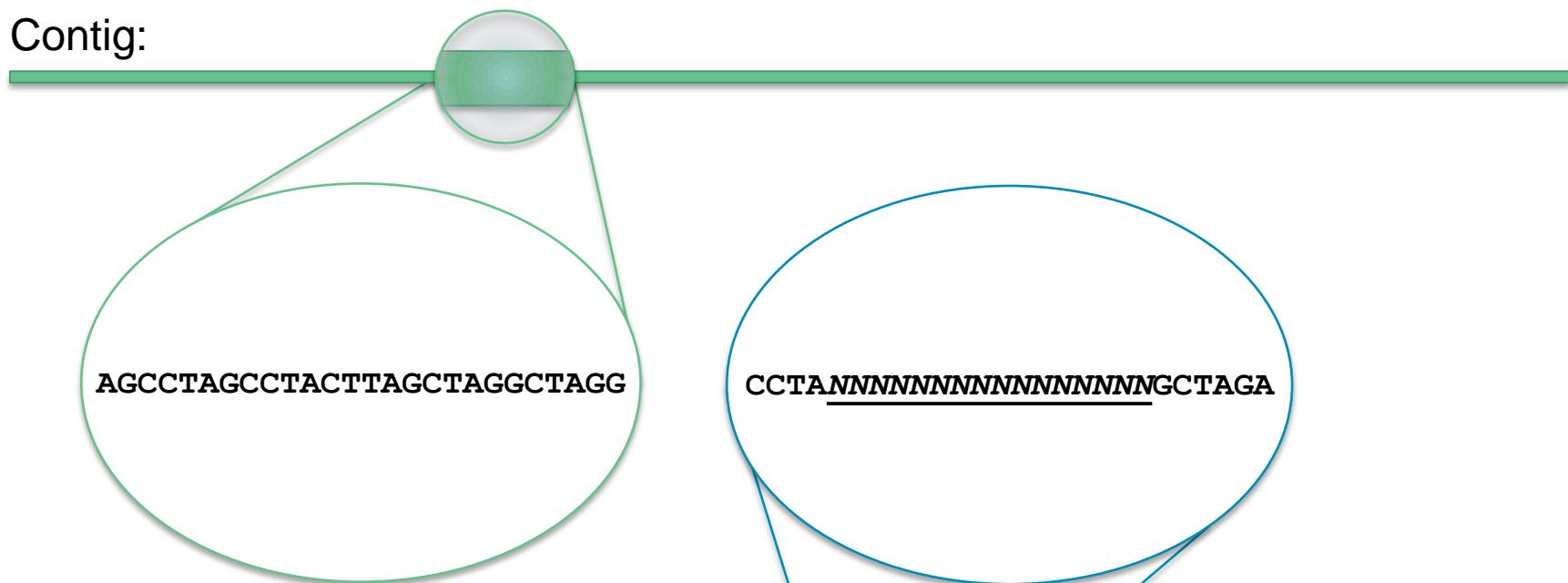
¹ Kawada et al. (2016) Discovery of Point Mutations in the Voltage-Gated Sodium Channel from African *Aedes aegypti* Populations: Potential Phylogenetic Reasons for Gene Introgression. PLoS Negl Trop Dis 10: e0004780

² ftp://ftp.ensemblgenomes.org/pub/metazoa/release-31/fasta/aedes_aegypti/dna/

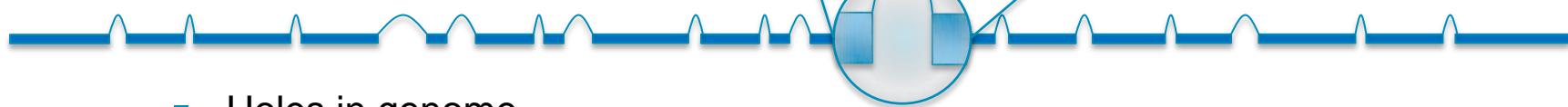
³ ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA_001014885.1_ASM101488v1

CONTIGS vs. SCAFFOLDS

Contig:



Scaffold:



- Holes in genome
 - Missing/incorrect information
 - Knowledge gaps

SUMMARY

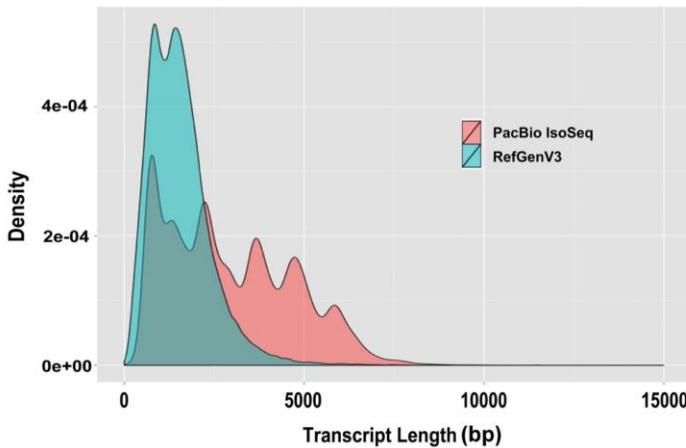
Advantages of Diploid Assembly with PacBio

- More accurate assembly of non-model or non-laboratory organisms
 - Non-inbred
 - Heterozygous
 - Pooled individuals
- Phasing of long-range allelic haplotypes
- More complete base-pair resolution of genome
- Higher genome contiguity and completeness
- Estimate of variation in parental chromosomes
 - Spatial patterns of variation across the genome
 - Copy number/structural variation

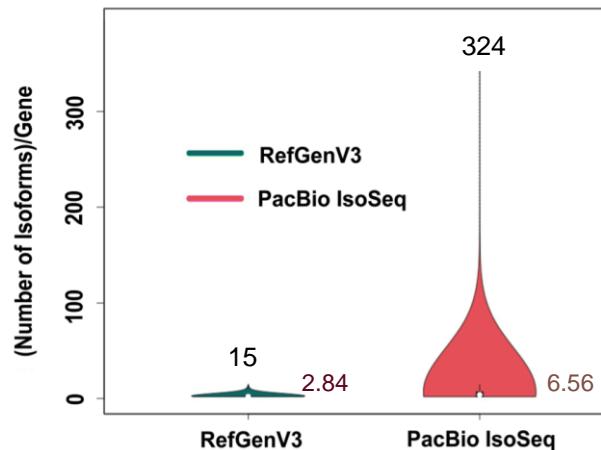
GENOME ANNOTATION

Iso-Seq Method for Full-Length *de novo* Transcript Assembly

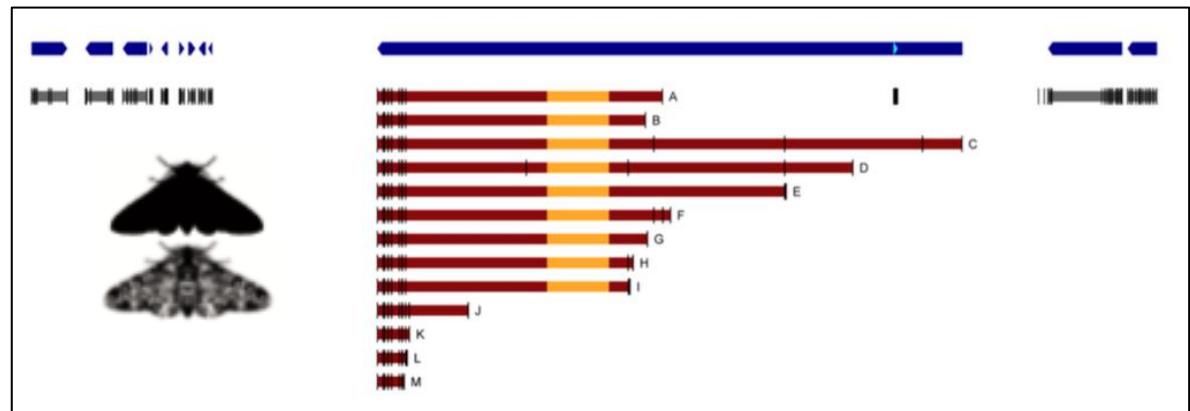
- Longer transcripts



- More numerous isoforms

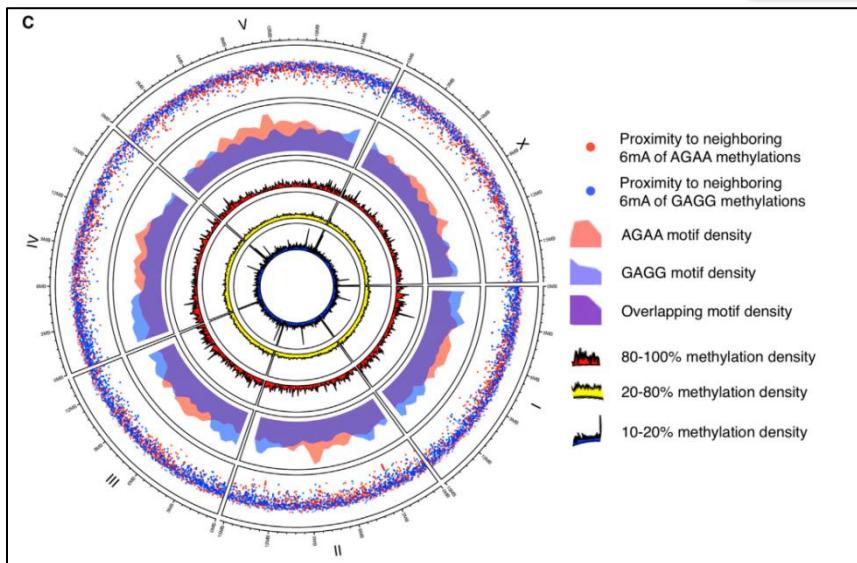
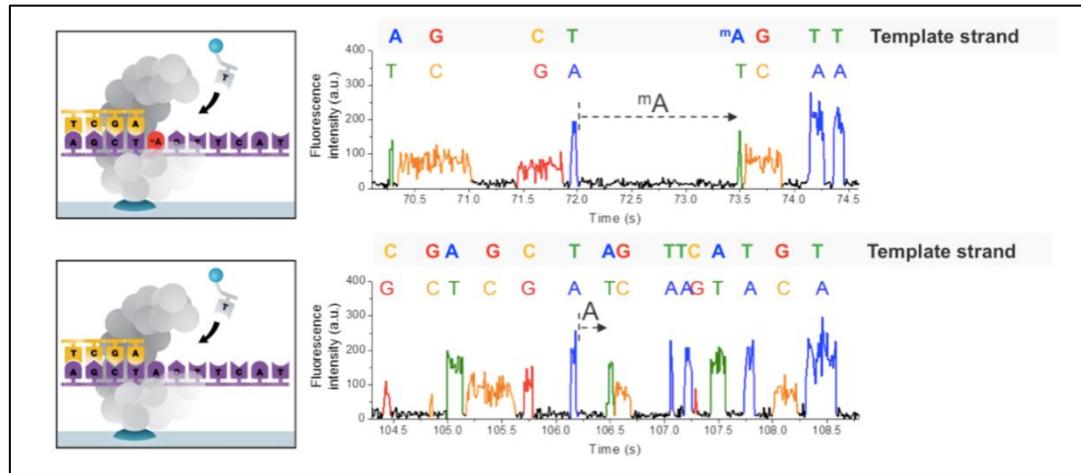


- Peppered moth
"carbonaria" morph due to TE insertion in *cortex* gene



EPIGENOME CHARACTERIZATION

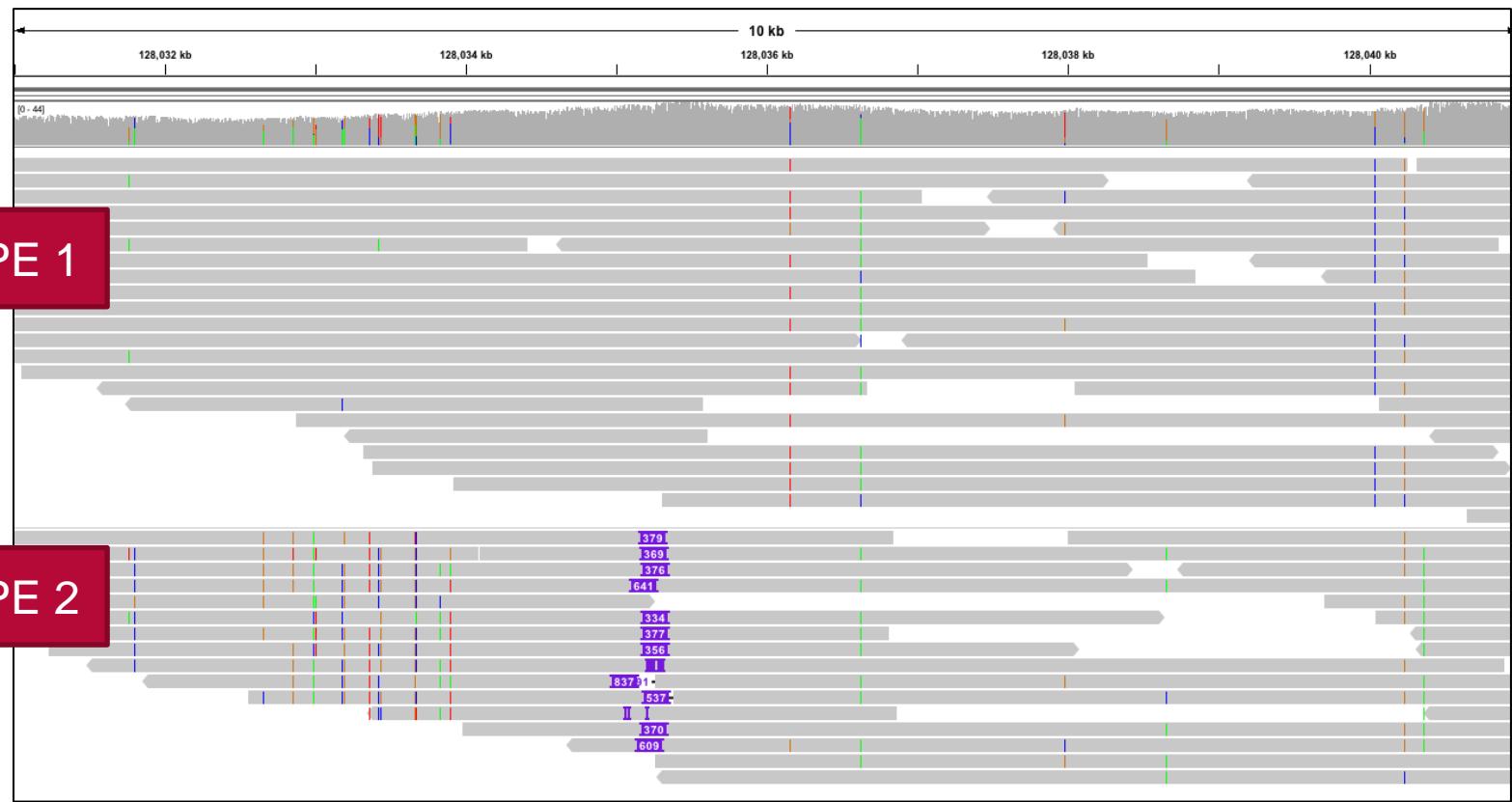
- Detection of base modifications
- Common mods: 5mC, 4mC, 6mA
- Detected through polymerase kinetics



- Relationship to gene expression, host-pathogen interactions, sequence motifs, DNA-damage/repair

TOOLS: IGV FOR VARIANT VISUALIZATION

Single Nucleotide Variants and Haplotype Resolution



TOOLS: IGV FOR VARIANT VISUALIZATION

Single Nucleotide Variants and Haplotype Resolution

HAPLOTYPE 1

HAPLOTYPE 2

Development Snapshot Build *Latest development snapshot; built at least nightly.*

Source code repository is hosted at GitHub. Current release branch is 2.3.x:

- <https://github.com/igvteam/igv/>

License

Permission to use this work is granted under the [MIT License](#)

Contact

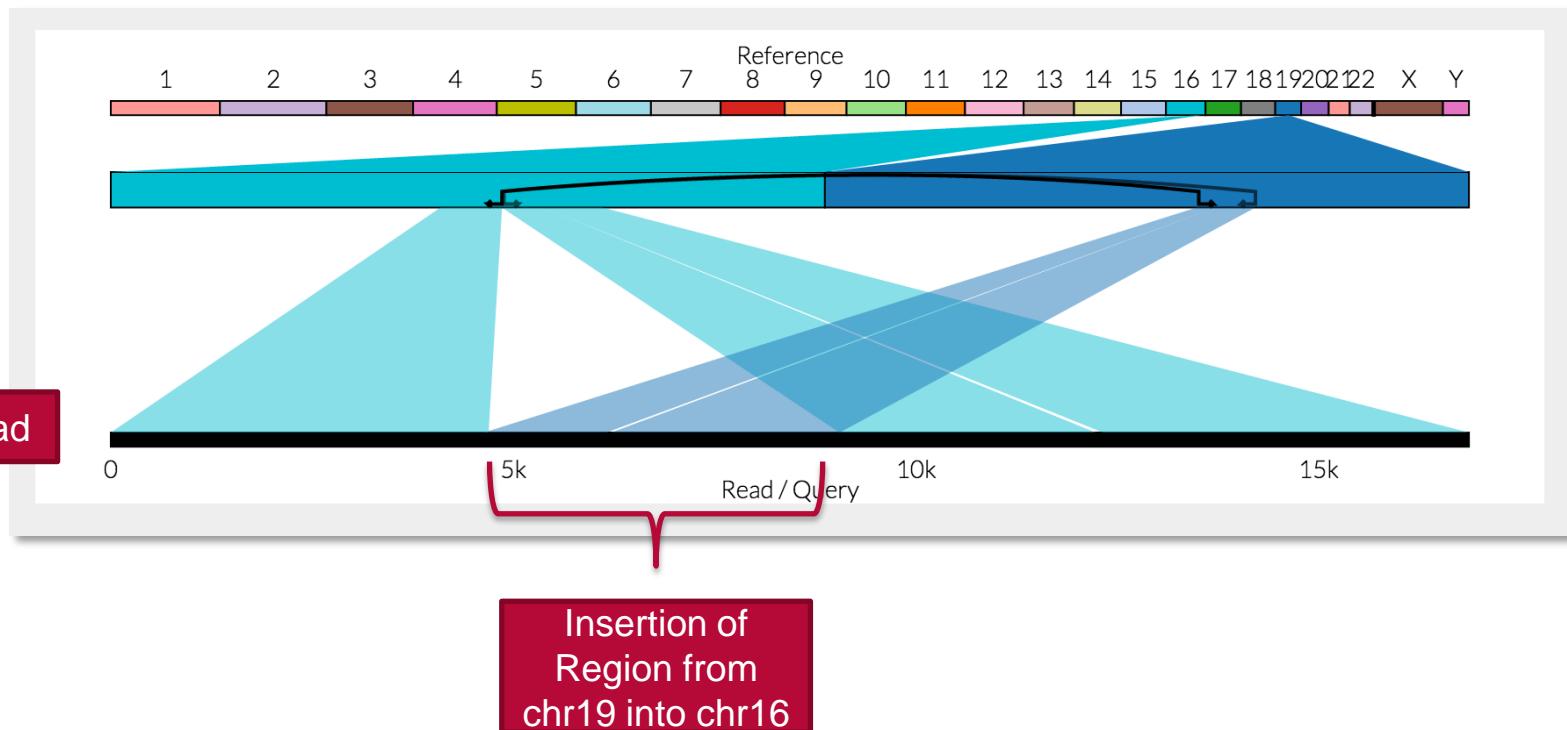
We welcome your feedback. Please send suggestions, requests, and bug reports to our [help forum](#) or open a ticket in our [issue tracker](#).

If you are posting about a problem, please include the error message (if applicable) and your `igv.log` file. It will be in <your home directory>/igv. When posting a bug report, please include:

1. The version of IGV you are using
2. What you did (e.g. I opened a bedgraph file). If the problem happens when you open a file and/or view certain genomic coordinates, attach the file* (or provide a link) and specify the genomic coordinates.
3. What you expected to happen (e.g. I expected IGV to display my data)
4. What actually happened (e.g. I received an error message saying "Unknown extension: .graph")

GENOME ANNOTATION

Structural Variation Visualized in Ribbon



SUMMARY

Advantages of Diploid Assembly with PacBio

- More accurate assembly of non-model or non-laboratory organisms
 - Non-inbred
 - Heterozygous
 - Pooled individuals
- Phasing of long-range allelic haplotypes
- More complete base-pair resolution of genome
- Higher genome contiguity and completeness
- Estimate of variation in parental chromosomes
 - Spatial patterns of variation across the genome
 - Copy number/structural variation