



# Assembly and annotation of diploid and polyploid genomes with PacBio

Sarah B Kingan, Bioinformatics Scientist, PacBio Applications  
January 12<sup>th</sup>, 2018, San Diego Botanical Garden

---

# AGENDA

- Intro to PacBio data for genome assembly and annotation
- Assembly workflow using FALCON-Unzip
- Understanding assembly output for complex genomes

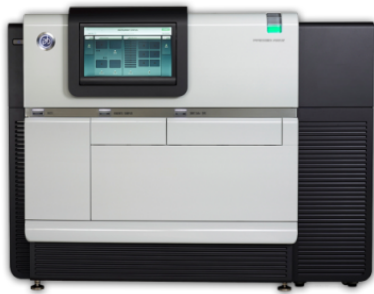
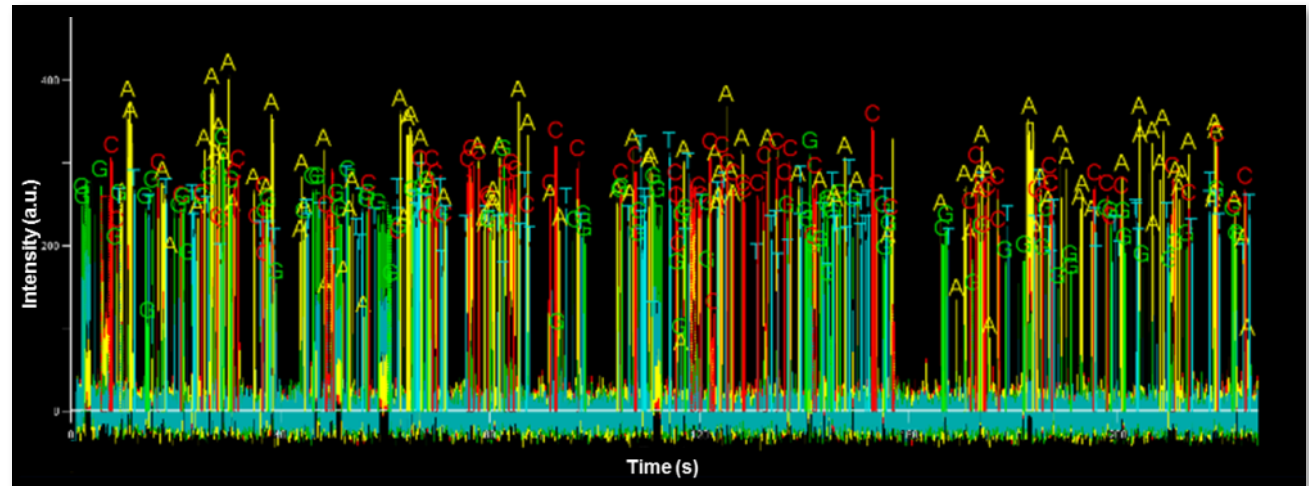
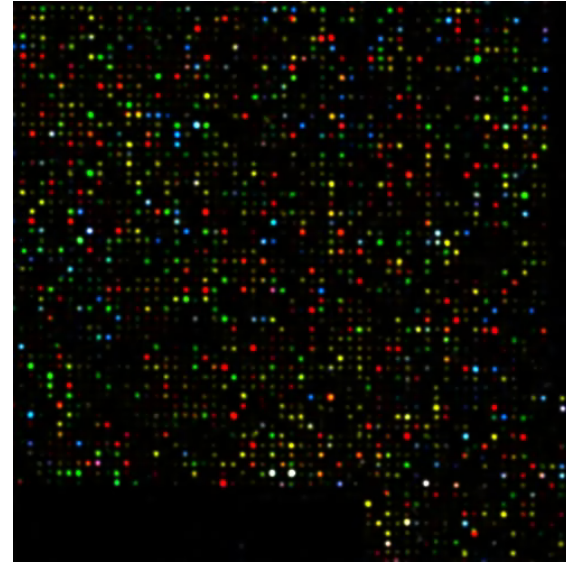
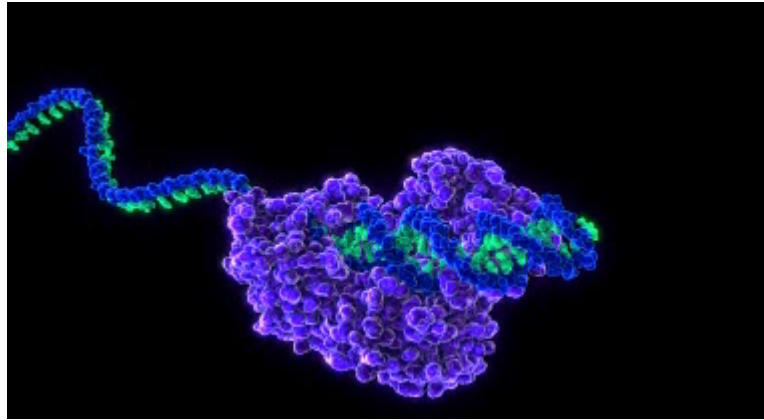


# Intro to PacBio data

Applications for genome assembly and annotation



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





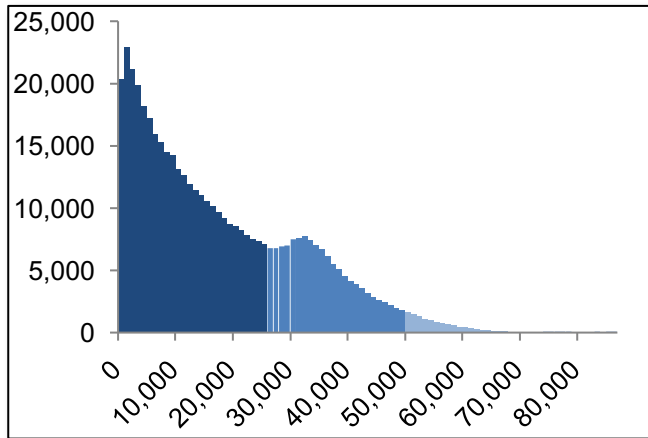
## SEQUEL SYSTEM

- Throughput per Cell: 5 – 10 Gb
- Average read length: 10 – 20 kb
- Read per cell: 400,000
- SMRT Cells per run: 1 – 16
  
- Improved performance with new chemistry and software release February 2018



# SMRT SEQUENCING CHARACTERISTICS

## Read Length Histogram

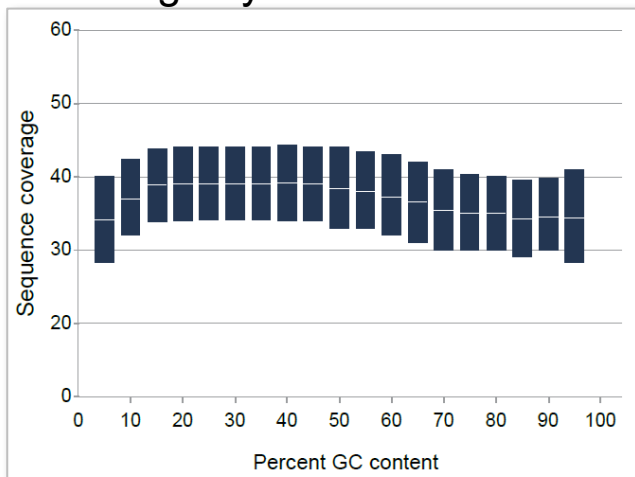


Rice 30 kb size-selected library using the Express kit, Sequel System with 2.1 Chemistry, 5.1 Sequel System Software.

## Long Reads

- Resolve repeats
- Contiguous, gapless contig assemblies
- Long-range haplotype phasing

## Coverage by GC%



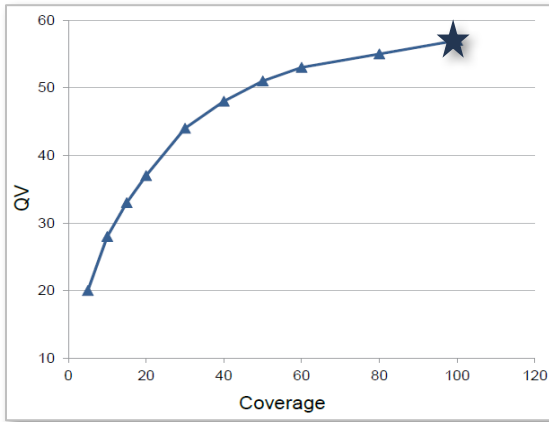
40 kb human library on a Sequel System using 2.1 chemistry and SMRT® Analysis v 5.1

## Uniform, Unbiased Coverage

- Sequence *entire* genome
- Longer, more complete assemblies

# SMRT SEQUENCING CHARACTERISTICS

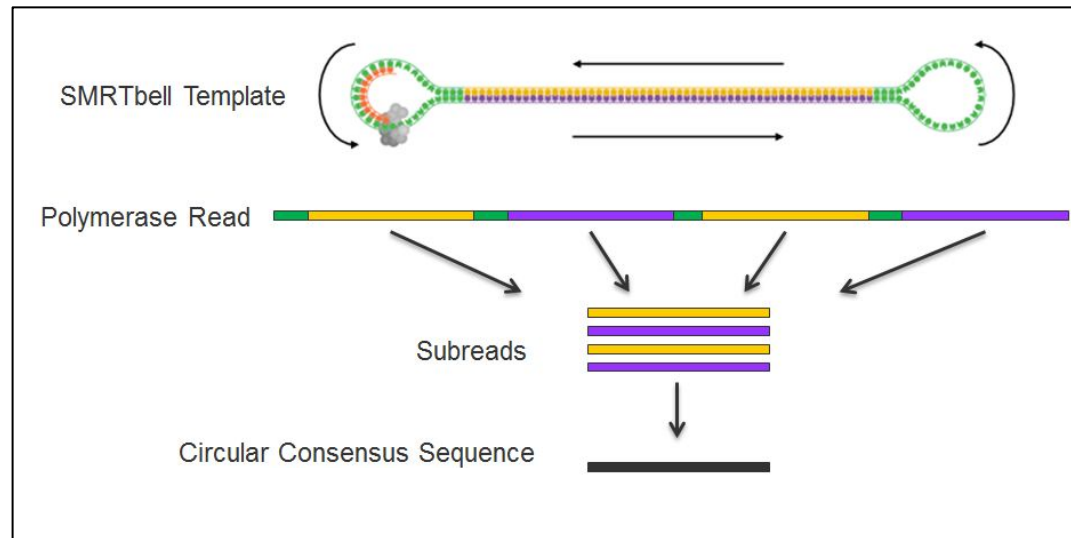
## QV by Coverage



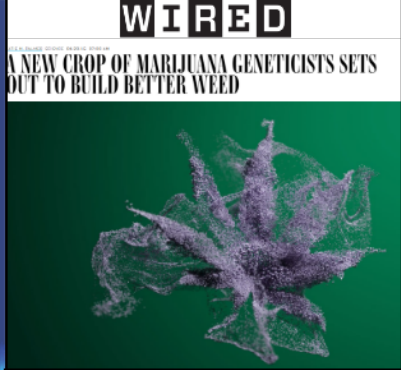
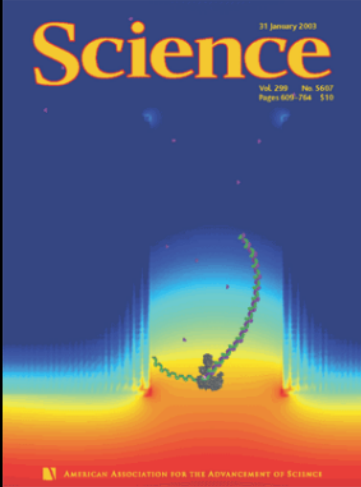
25 kb *E. coli* library on a Sequel System using 2.1 chemistry and SMRT® Analysis v 5.1

## High Consensus Accuracy

- Random error profile
- Achieves QV50
- 99.999% accuracy

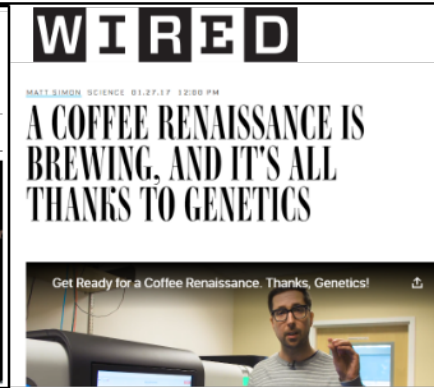
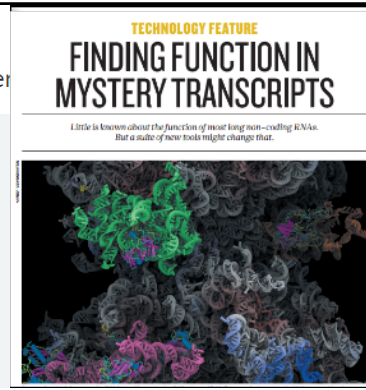






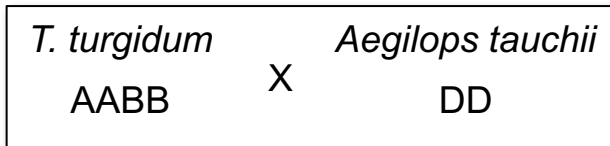
## PACBIO APPLICATIONS

- Whole Genome Sequencing
- Isoform Sequencing (Iso-Seq Analysis)
- Structural Variants
- Targeted Sequencing
- Microbial Epigenetics



# WHY PACBIO FOR *DE NOVO* GENOME ASSEMBLY?

- *Triticum aestivum* (bread wheat)
- Genome size >15 Gb
- allohexaploid (AABBDD)



Assembly	Data	Length	Contig N50
IWGSC <sup>1</sup>	100-fold ILM	10.2 Gb	8.9 kb
FALCON <sup>2</sup>	36-fold PB	12.9 Gb	215 kb
MaSuRCA <sup>2</sup>	36-fold PB + 64-fold ILM	17.0 Gb	76 kb
Merged <sup>2</sup>	NA	15.3 Gb	233 kb

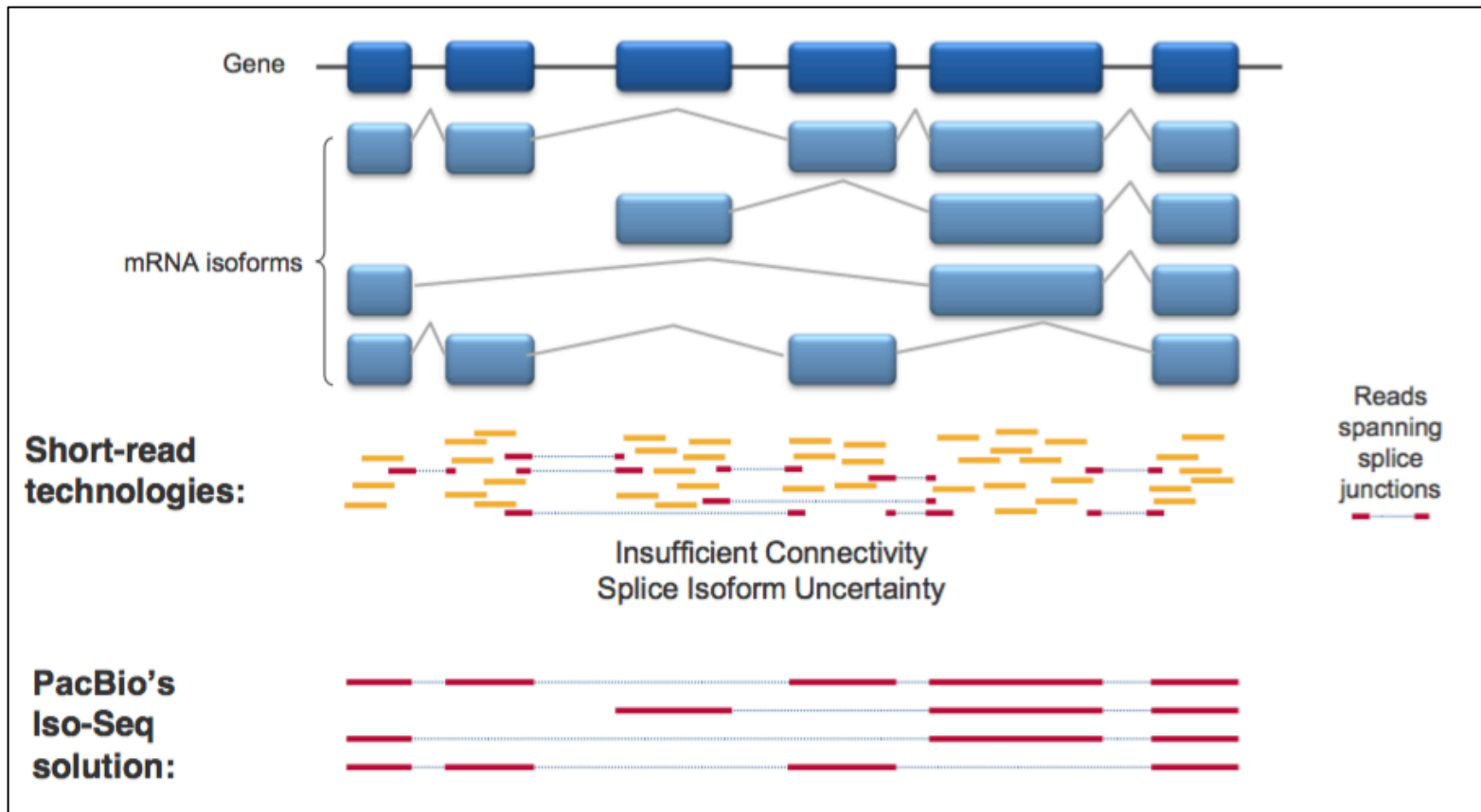


1. IWGSC (2014) Science 345:6194.  
 2. Zimin et al. (2017) GigaScience 6:1

# WHY PACBIO FOR GENOME ANNOTATION?

- Isoform Sequencing (Iso-Seq Analysis) aka RNA-seq
- Full Length cDNA sequences
- No assembly required

**>100 Publications using Iso-Seq Analysis**





# ISO-SEQ ANALYSIS FOR GENOME ANNOTATION

- Whole RNA extracted from brain
- 2 Sequel cells per sample
- ~400,000 full length isoforms

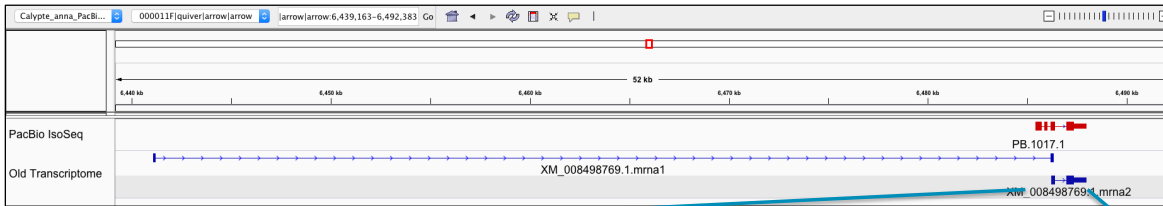


Zebra Finch  
*Taeniopygia guttata*

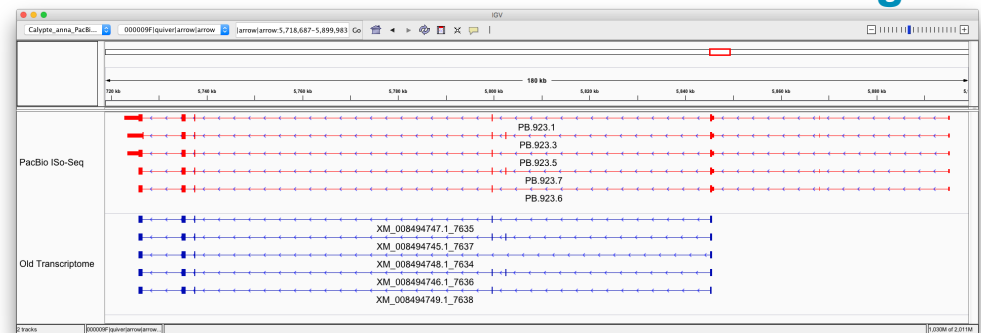


Anna's Hummingbird  
*Calypte anna*

## Corrected Gene Model: DUSP1



## Extended UTRs: neuroligin



# ISO-SEQ ANALYSIS FOR GENOME ANNOTATION

- Whole RNA extracted from brain
- 2 Sequel cells per sample
- ~400,000 full length isoforms

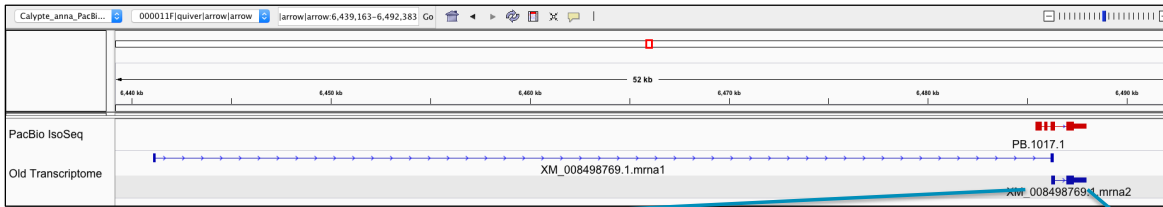


Zebra Finch  
*Taeniopygia guttata*



Anna's Hummingbird  
*Calypte anna*

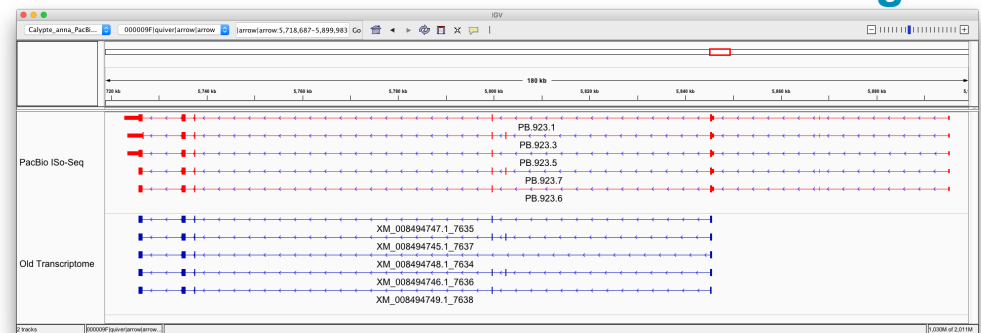
## Corrected Gene Model: DUSP1



1-2 cells per tissue for genome annotation

Analysis in SMRT Link GUI

## Extended UTRs: neuroigin



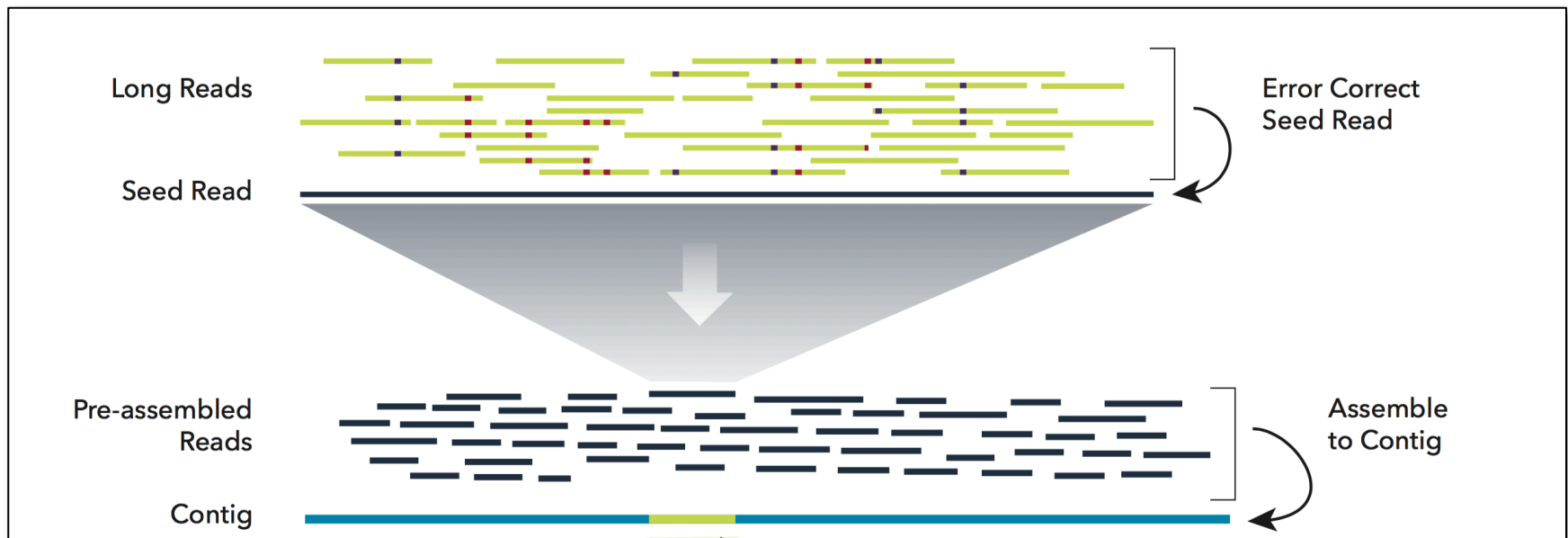
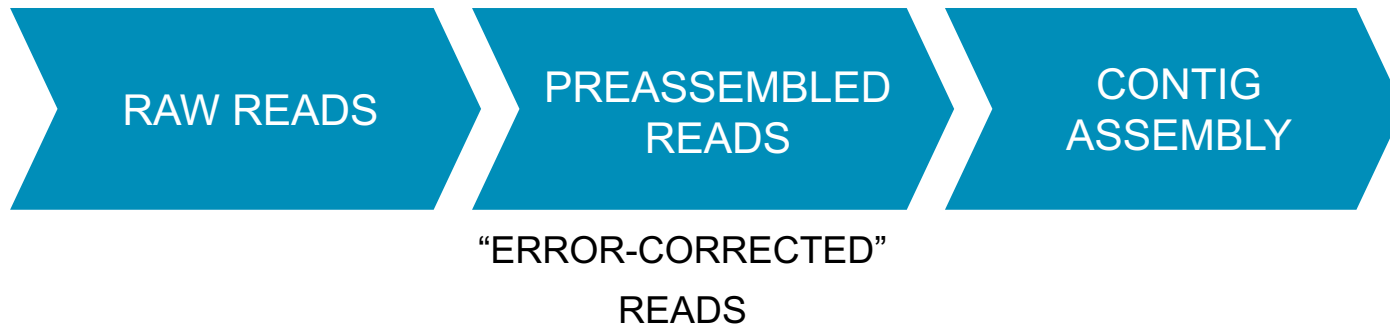


# *de novo* Assembly Workflow

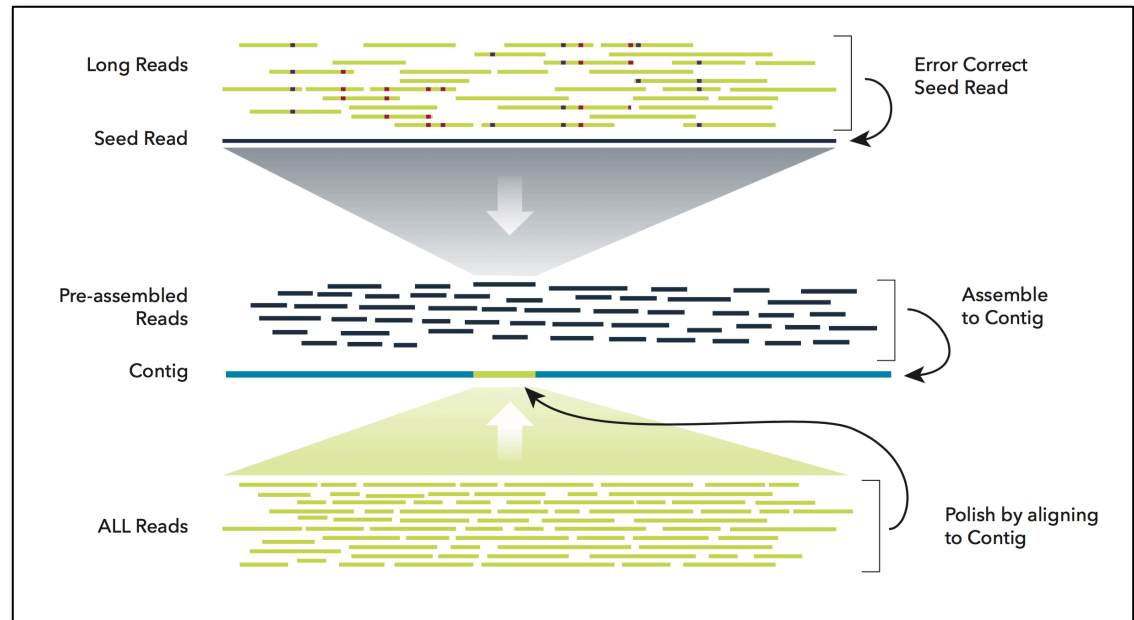
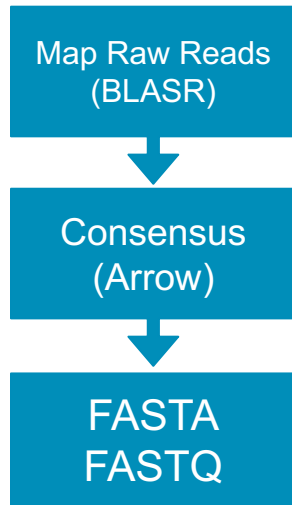
FALCON and FALCON-Unzip for phased contig assembly



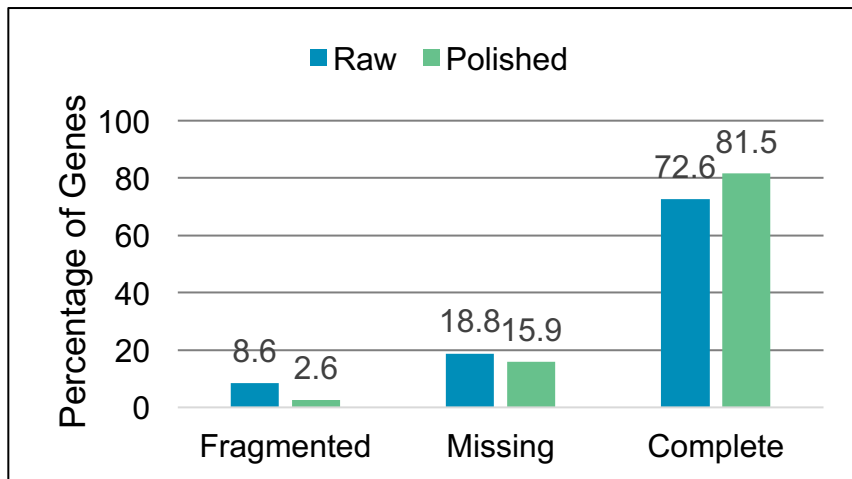
# FALCON / HIERARCHICAL GENOME ASSEMBLY PROCESS (HGAP)



# POLISHING WITH PACBIO DATA IMPROVES BASE QUALITY



## GENOME COMPLETENESS WITH BUSCO



70% reduction in **Fragmented Genes**  
 15% reduction in **Missing Genes**  
 12% increase in **Complete Genes**

**Acknowledgement:**  
 Erich Jarvis, Rockefeller University

# FALCON AND FALCON-UNZIP



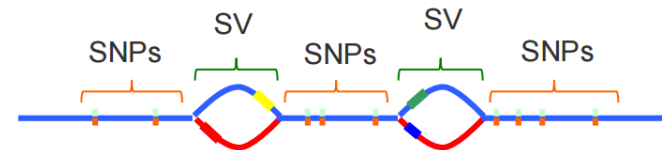
## Phased diploid genome assembly with single-molecule real-time sequencing.

Chin CS, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, Dunn C, O'Malley R, Figueroa-Balderas R, Morales-Cruz A, Cramer GR, Delledonne M, Luo C, Ecker JR, Cantu D, Rank DR, Schatz MC

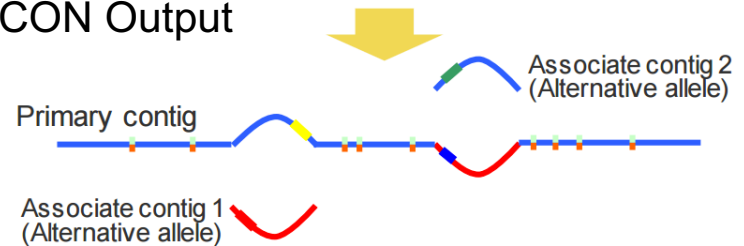
### ABSTRACT

While genome assembly projects have been successful in many haploid and inbred species, the assembly of noninbred or rearranged heterozygous genomes remains a major challenge. To address this challenge, we introduce the open-source FALCON and FALCON-Unzip algorithms (<https://github.com/PacificBiosciences/FALCON/>) to assemble long-read sequencing data into highly accurate, contiguous, and correctly phased diploid genomes. We generate new reference sequences for heterozygous samples including an F1 hybrid of *Arabidopsis thaliana*, the widely cultivated *Vitis vinifera* cv. Cabernet Sauvignon, and the coral fungus *Clavicornia pyxidata*, samples that have challenged short-read assembly approaches. The FALCON-based

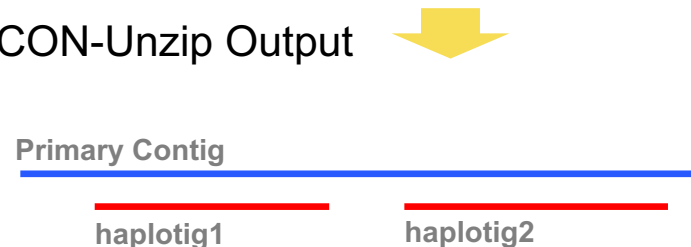
### Initial Assembly Graph



### FALCON Output



### FALCON-Unzip Output



- FALCON is a **diploid-aware assembler**.
- FALCON-Unzip module performs true **phased assembly** for diploid or polyploid samples.

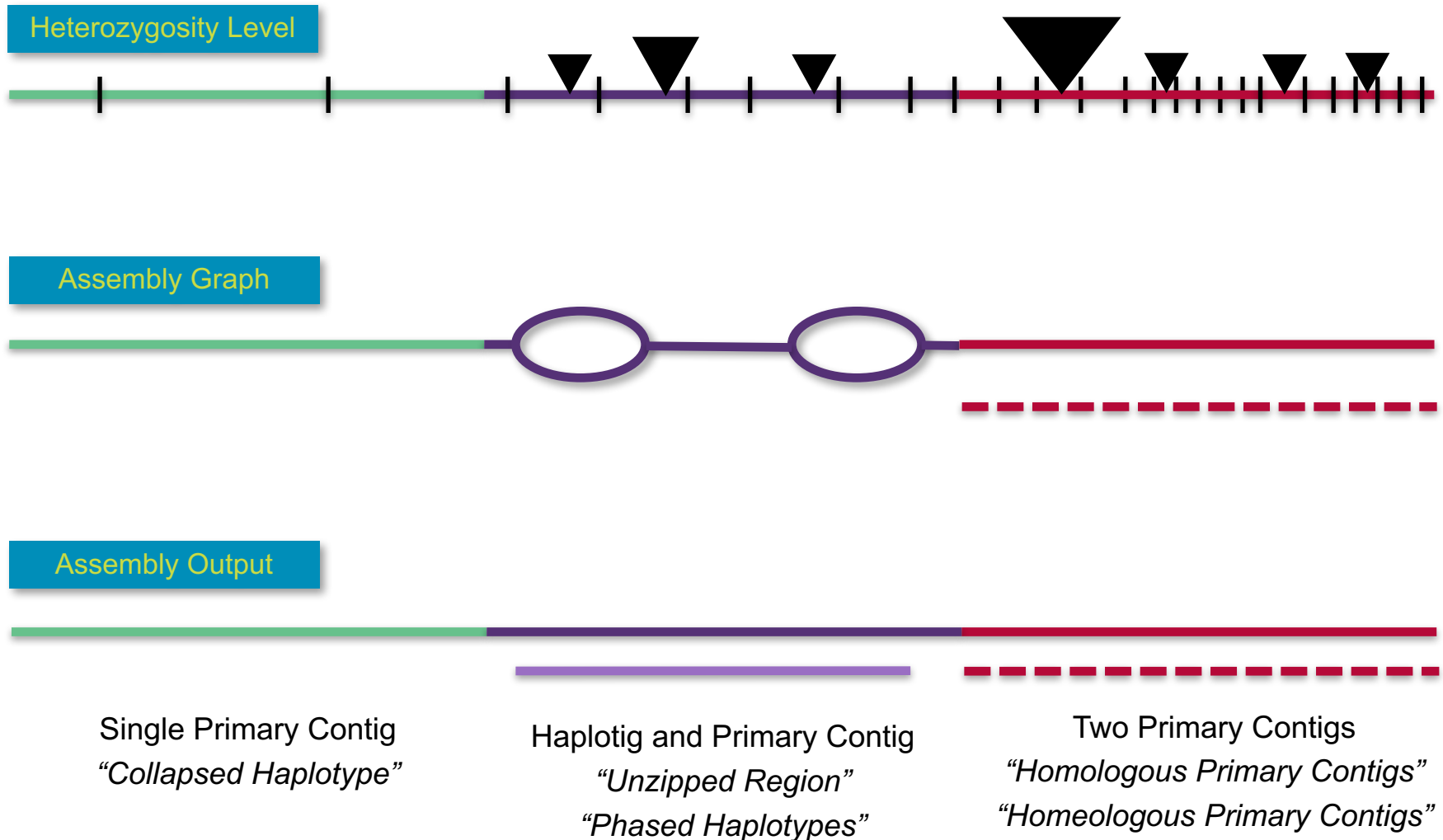




# Understanding Complex Assemblies

Leveraging coverage, gene annotation, and alignments

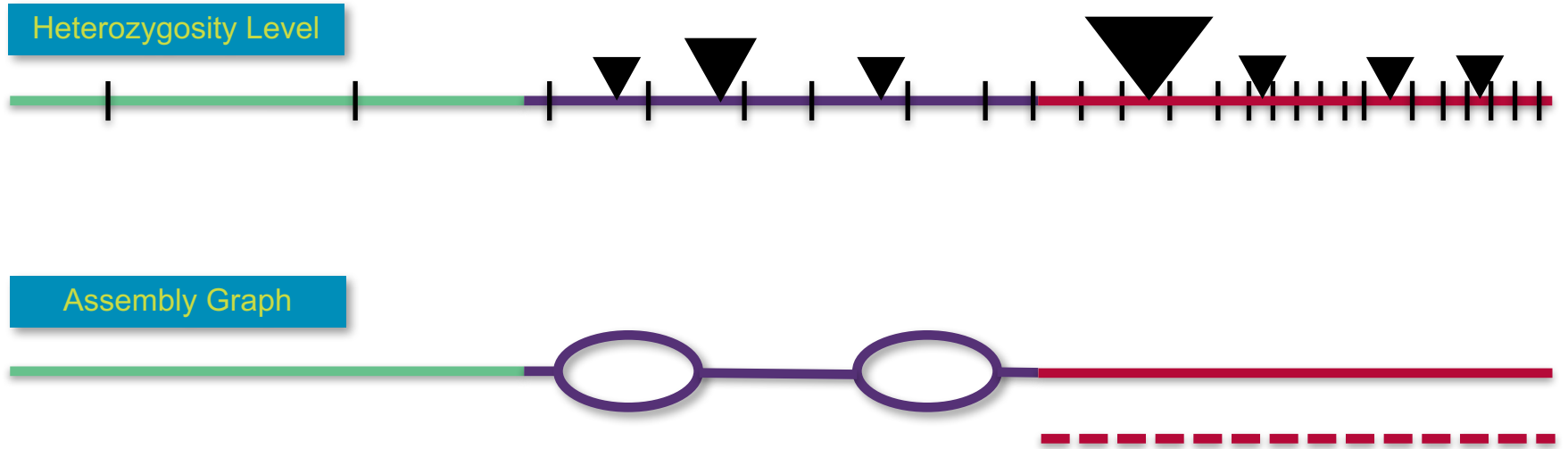
# IMPACT OF HETEROZYGOSITY ON ASSEMBLY PROCESS



# IMPACT OF HETEROZYGOSITY ON ASSEMBLY PROCESS

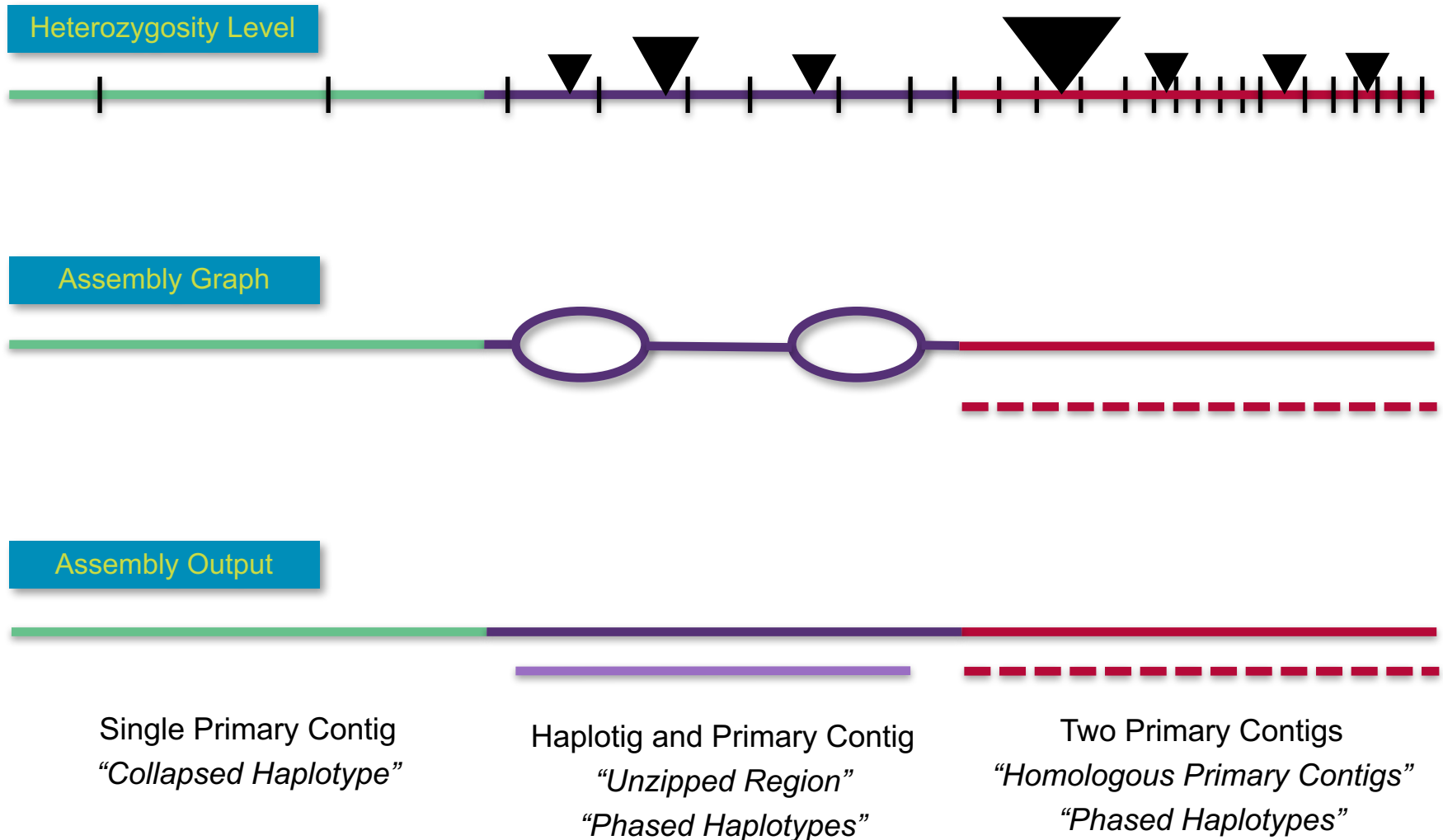


# IMPACT OF HETEROZYGOSITY ON ASSEMBLY PROCESS

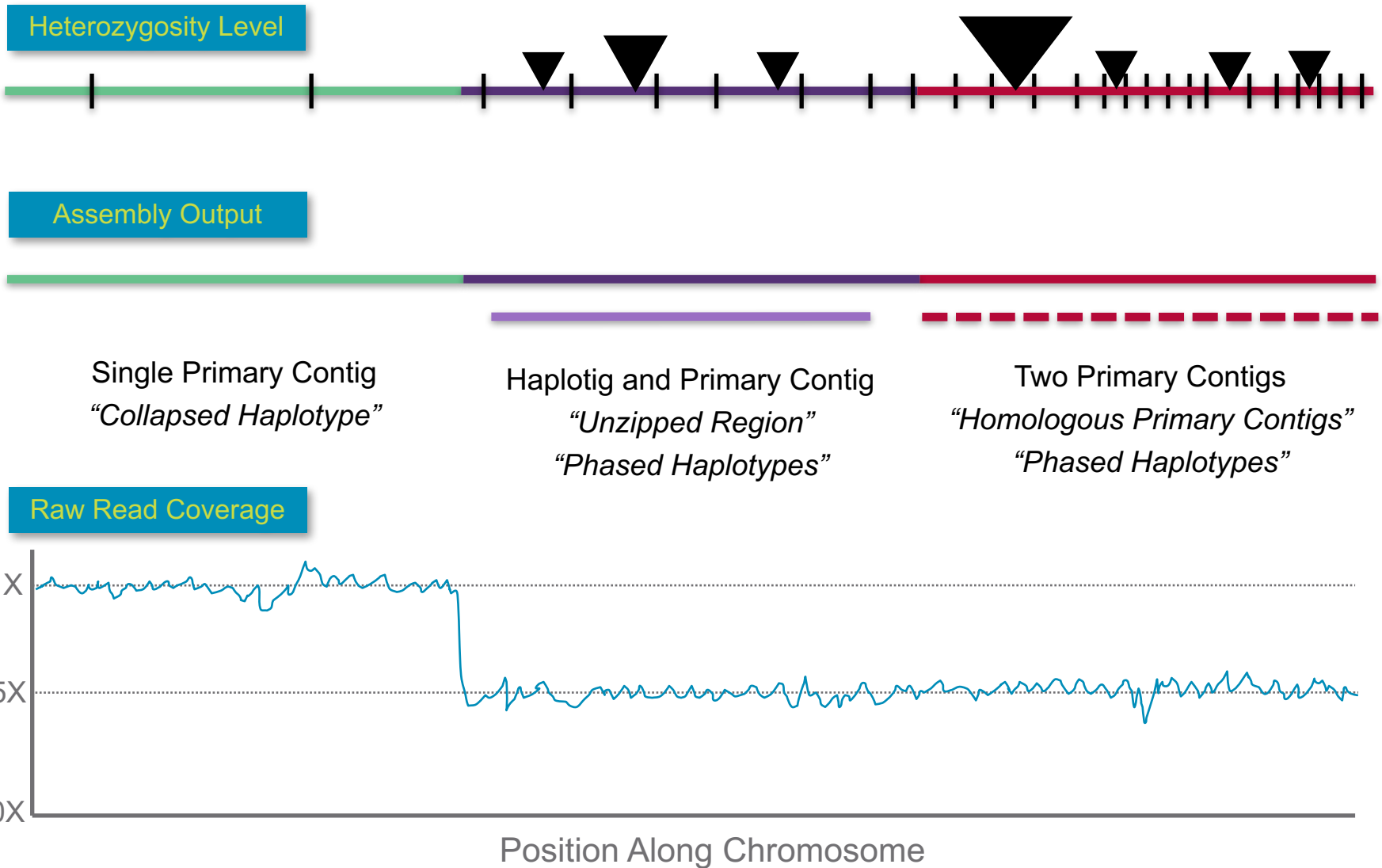




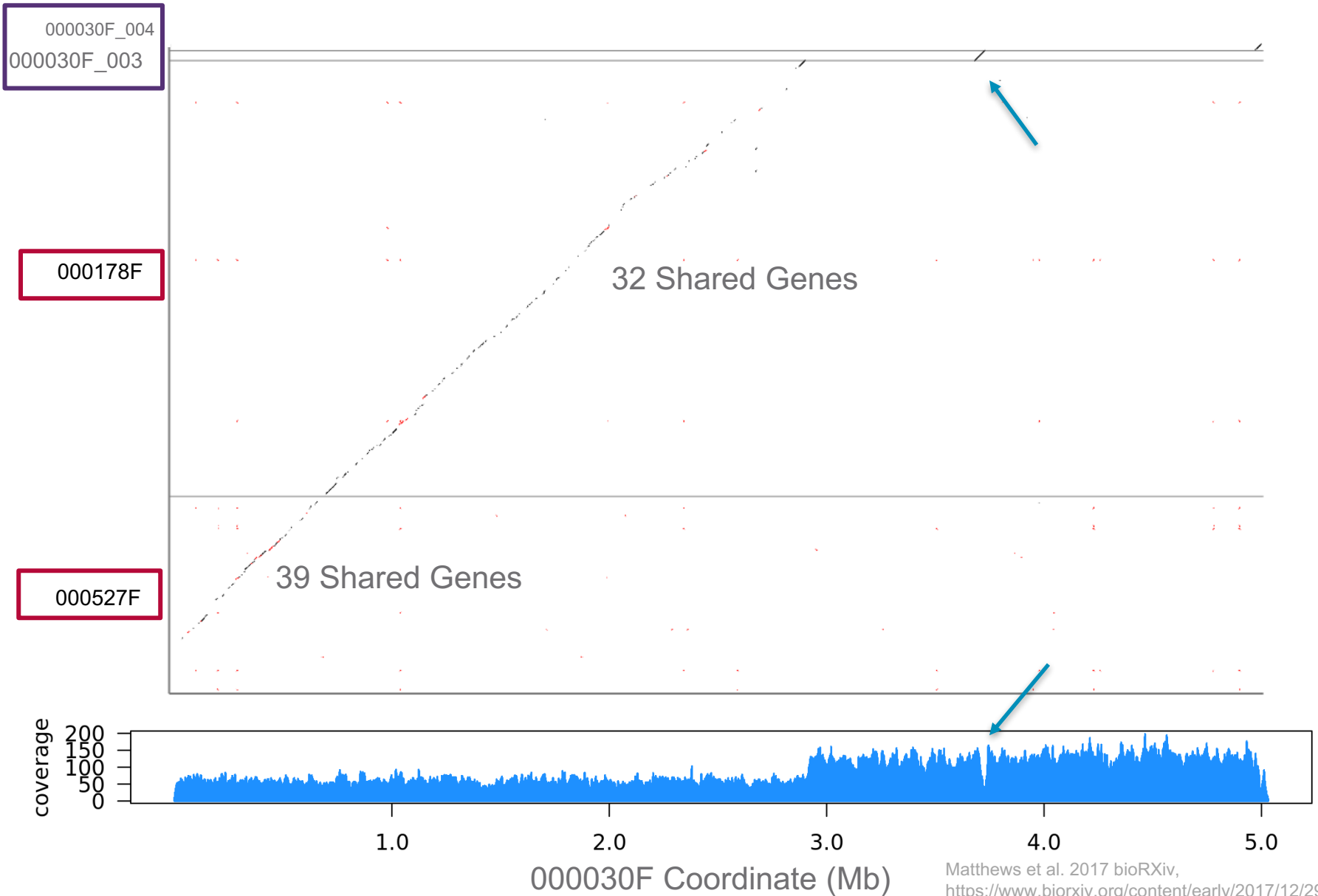
# IMPACT OF HETEROZYGOSITY ON ASSEMBLY PROCESS



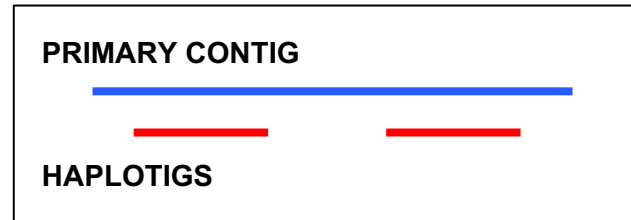
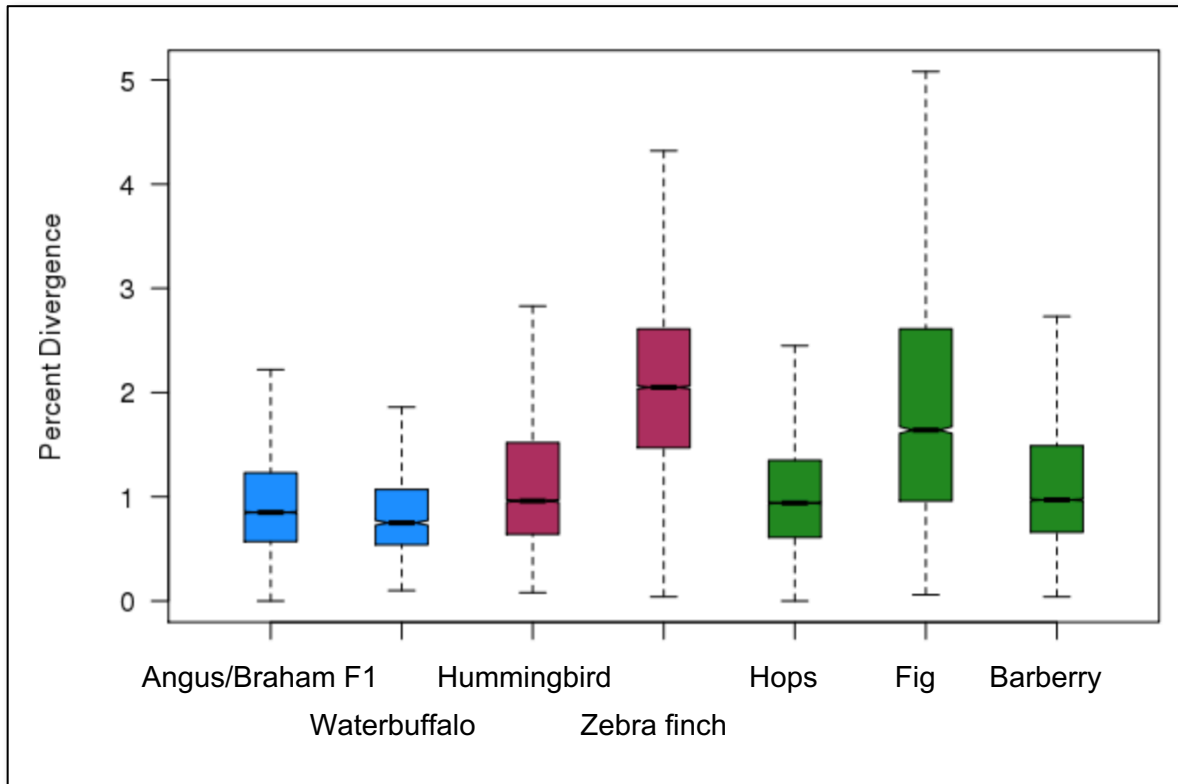
# RAW READ COVERAGE AND ASSEMBLY STRUCTURE



# HOMOLOGOUS PRIMARY CONTIGS IN AEADES MOSQUITO



# HOW MUCH DIVERGENCE IS CAPTURED BY UNZIP?



- Up to 5% divergence captured in “Unzipped” regions
- More divergent haplotypes assembled on separate primary contigs





[www.pacb.com](http://www.pacb.com)

For Research Use Only. Not for use in diagnostics procedures. © Copyright 2018 by Pacific Biosciences of California, Inc. All rights reserved. Pacific Biosciences, the Pacific Biosciences logo, PacBio, SMRT, SMRTbell, Iso-Seq, and Sequel are trademarks of Pacific Biosciences. BluePippin and SageELF are trademarks of Sage Science. NGS-go and NGSengine are trademarks of GenDx. FEMTO Pulse and Fragment Analyzer are trademarks of Advanced Analytical Technologies.

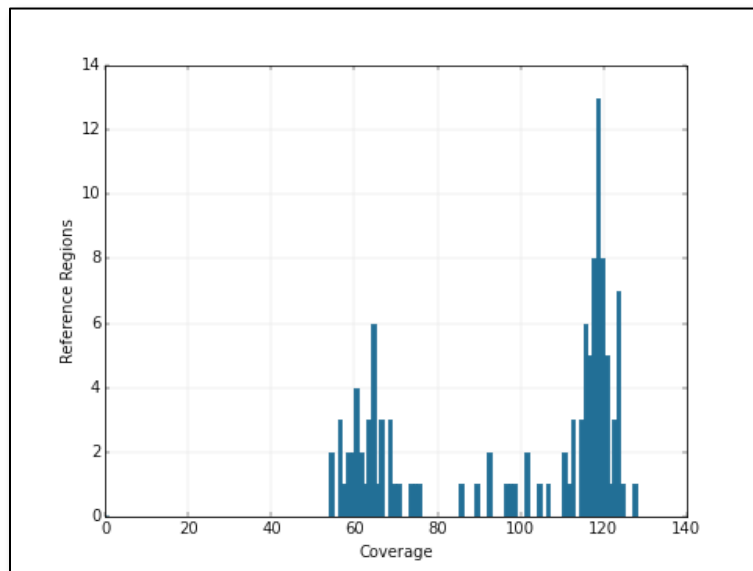
All other trademarks are the sole property of their respective owners.



# EXAMPLE: *Aedes* MOSQUITO FALCON-UNZIP ASSEMBLY

- Expected Genome Size: ~1.3 Gb
- Primary Contig Length: 1.69 Gb

## BIMODAL COVERAGE HISTOGRAM



## BUSCO ANALYSIS:

### ARTHROPOD GENESET (N = 2675)

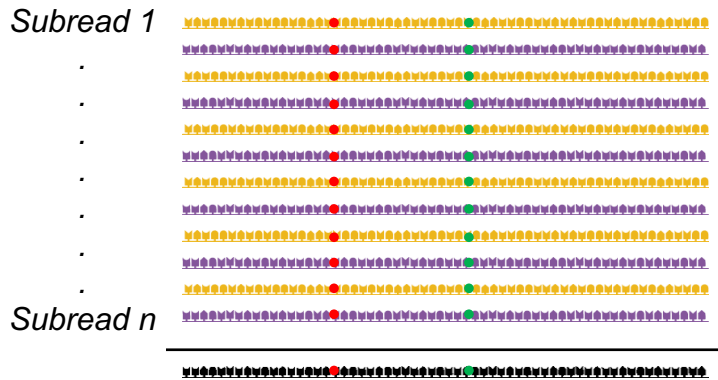
ASSEMBLY	<i>Aedes</i> PACBIO
COMPLETE	98%
MISSING	2%
FRAGMENTED	10%
DUPLICATED	32%

### Acknowledgement:

***Aedes* Genome Working Group  
 Leslie Vosshall, Ben Matthews,  
 Rockefeller University**

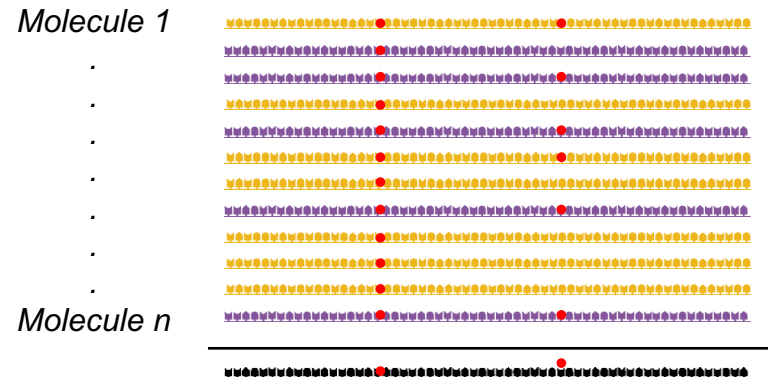
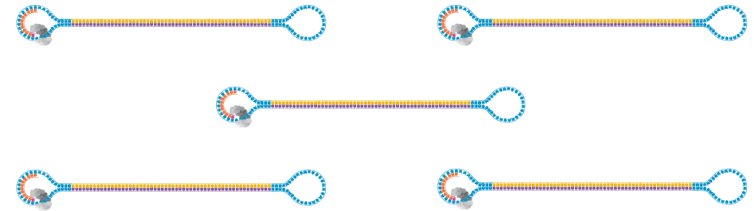
# MULTI- vs SINGLE-MOLECULE CONSENSUS

— Circular consensus sequencing (CCS):



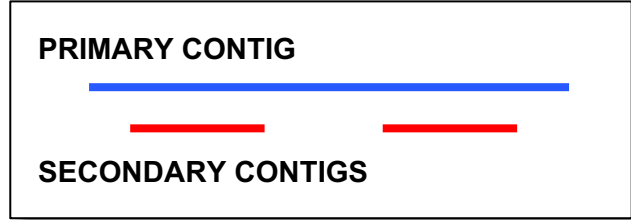
RNA-seq/Iso-Seq,  
targeted sequencing

— Large insert sequencing:



*de novo* assembly,  
SV detection

# EXAMPLE ASSEMBLY OF WATER BUFFALO



	FALCON-Unzip	FALCON	Williams et al. 2017 <sup>1</sup>
Primary Length	2.65 Gb	2.66 Gb	2.09 Gb
Primary N50	18.8 Mb	18.7 Mb	0.022 Mb
Secondary Length	1.53 Gb	0.218 Gb	NA
Proportion Phased	58 %	8.2 %	NA



**Olimpia**

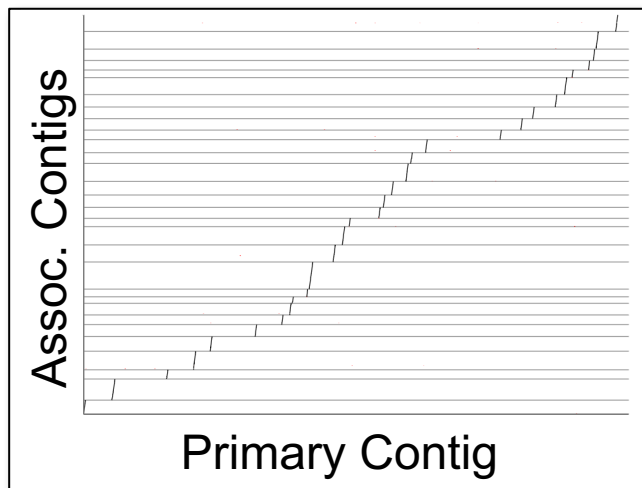
Photo Credit: Caterina Cambuli

**7-fold increase in haplotype phasing with Unzip module**

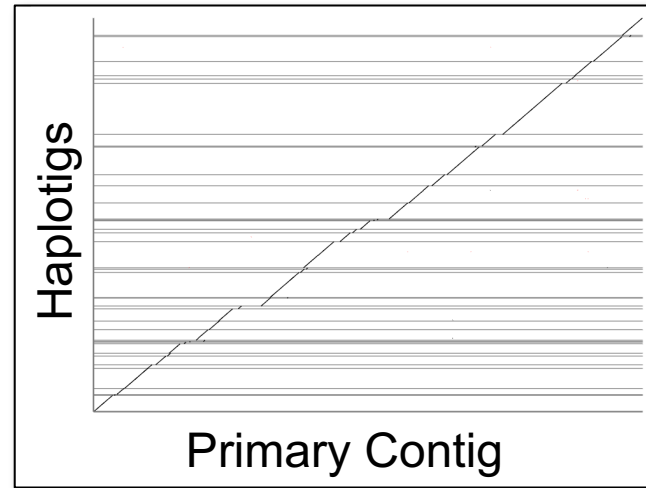
Acknowledgements:  
 Tim Smith, USDA-ARS  
 John Williams, Lloyd Low, University of Adelaide  
 Paolo Ajmone-Marsan, Università Cattolica del S. Cuore  
 David Hume, Mick Watson, Roslin Institute  
 1. Williams et al (2017) Gigascience. 6(10):

# INCREASED HAPLOTIG CONTIGUITY WITH FALCON-UNZIP

**FALCON**



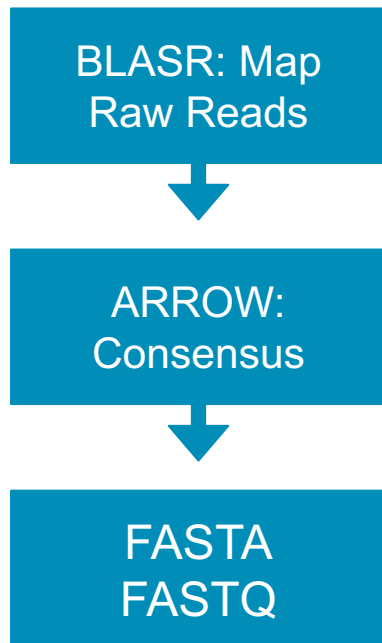
**FALCON-UNZIP**



<b>CONTIG: 000078F</b>	<b>FALCON</b>	<b>FALCON-Unzip</b>
<b>Primary Contig Length</b>	12.9 Mb	12.9 Mb
<b>Number Secondary Contigs</b>	30	34
<b>Total Secondary Length</b>	1.21 Mb	10.6 Mb
<b>Secondary Contig N50</b>	42.5 kb	470 kb
<b>Proportion Phased</b>	9.3 %	82%

# POLISHING WITH ARROW: WORKFLOW

METHOD	ASSEMBLY	POLISHING
HGAP4 - SMRT Link	✓	✓
FALCON	✓	resequencing pipeline from pbsmrtpipe/SMRT Link
FALCON-Unzip	✓	✓ (phased) plus optional resequencing



## Random Best Mapping

- Random choice of locus with equal BLASR score

## Minimum Coverage <5

- <5 reads span 500 bp window
- No consensus call
- Reference base returned as lowercase

Consensus Sequence

```

atgcgccggttatatgg
aagctagcTAGCTCTA
GTAGCTAGAGCTAGCT
GCGCGCTAGAA TAGGG
CGCCATAGAGCCTTTT
  
```



## ASSEMBLY METHOD RECOMMENDATIONS

METHOD	GENOME SIZE	HETEROZYGOSITY	COVERAGE
<b>HGAP4 - SMRT Link</b>	<3 GB*	Low	40-50 fold
<b>FALCON</b>	Any Size	Low - Medium	40-80 fold
<b>FALCON-Unzip**</b>	Any Size	Medium – High	80-100 fold
<b>Arrow Polishing</b>	ALWAYS POLISH 1-2 TIMES***		

\* Genome size limit depends on underlying compute resources.

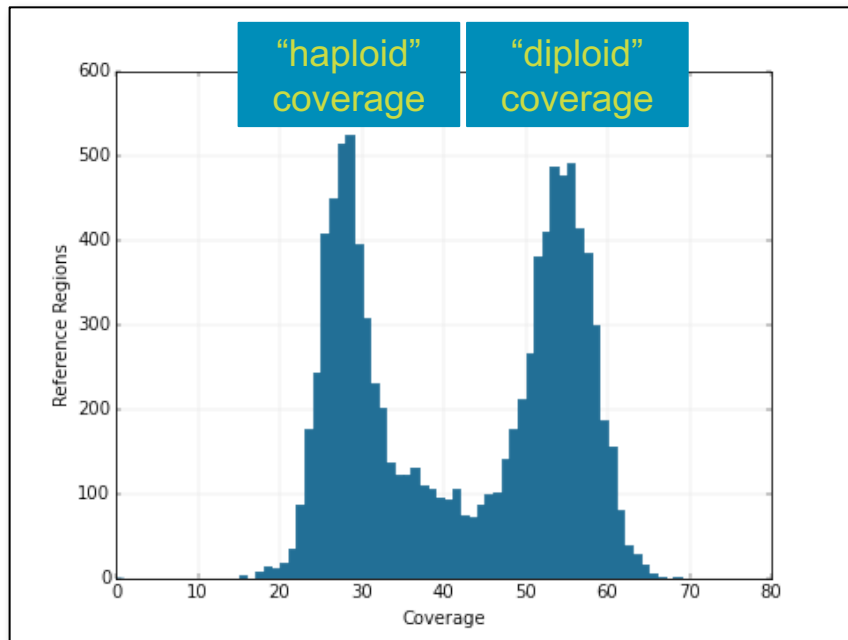
\*\* FALCON-Unzip must be run in a FALCON job directory. You CANNOT run HGAP4 and then FALCON-Unzip.

\*\*\* Reference sequence should be concatenated primary contigs and haplotigs

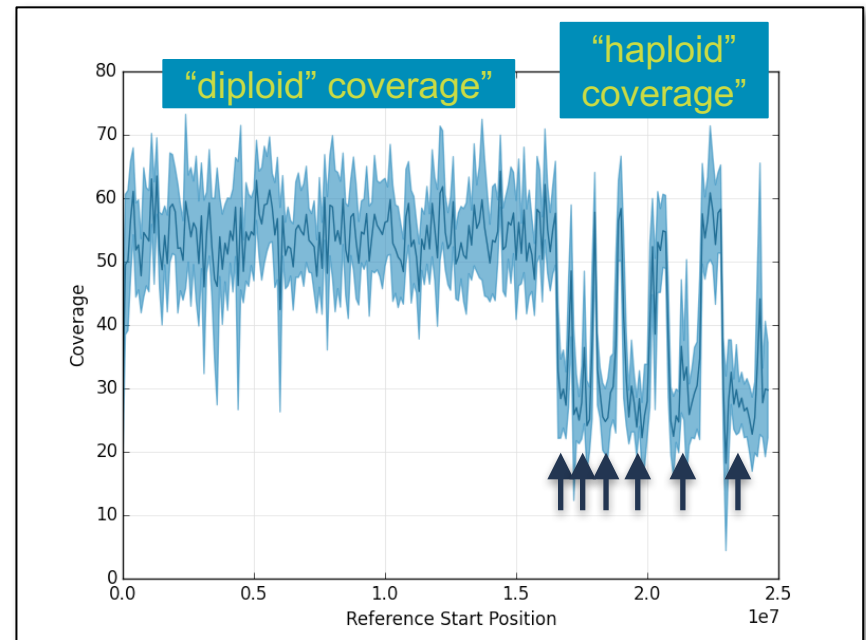
# SMRT LINK COVERAGE REPORTS

Graphical Outputs from Resequencing Pipeline / HGAP4

## COVERAGE HISTOGRAM: GENOME

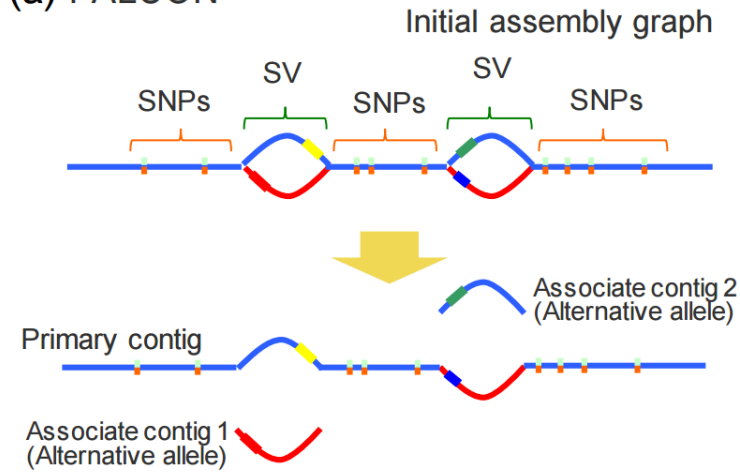


## COVERAGE PLOT: CONTIG

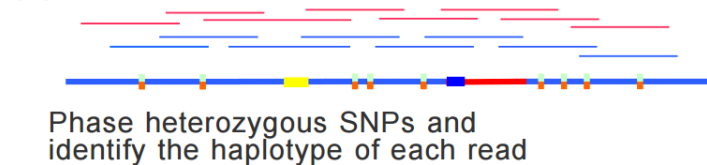


# DIPLOID ASSEMBLY WITH FALCON-UNZIP

(a) FALCON



(b)



(c) FALCON-Unzip

