

Using SMRT Iso-Seq Sequencing to Dissect Polyploid Transcriptomes: Lessons Learned from Tetra- and Hexaploid Blueberries



Hamid Ashrafi, PhD
Hamed Bostan, PhD
North Carolina State University
PacBio East Coast UGM - June 2017

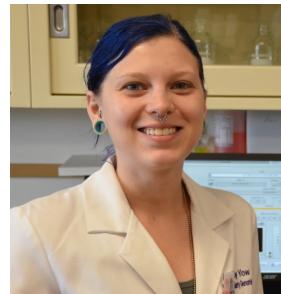
<https://blueberry.cals.ncsu.edu/>

<https://www.facebook.com/TeamVaccinium/>

Acknowledgements



Dr. George Yuan
PACIFIC BIOSCIENCES®



Ms. Ashley Yow
NCSU



Dr. Massimo Iorizzo
NCSU



Dr. Liz Tseng
PACIFIC BIOSCIENCES®



Dr. Rishi Aryal
NCSU



Dr. Hamed Bostan
NCSU

Funding Agencies



NC Blueberry Council NC Dept. Ag & CS

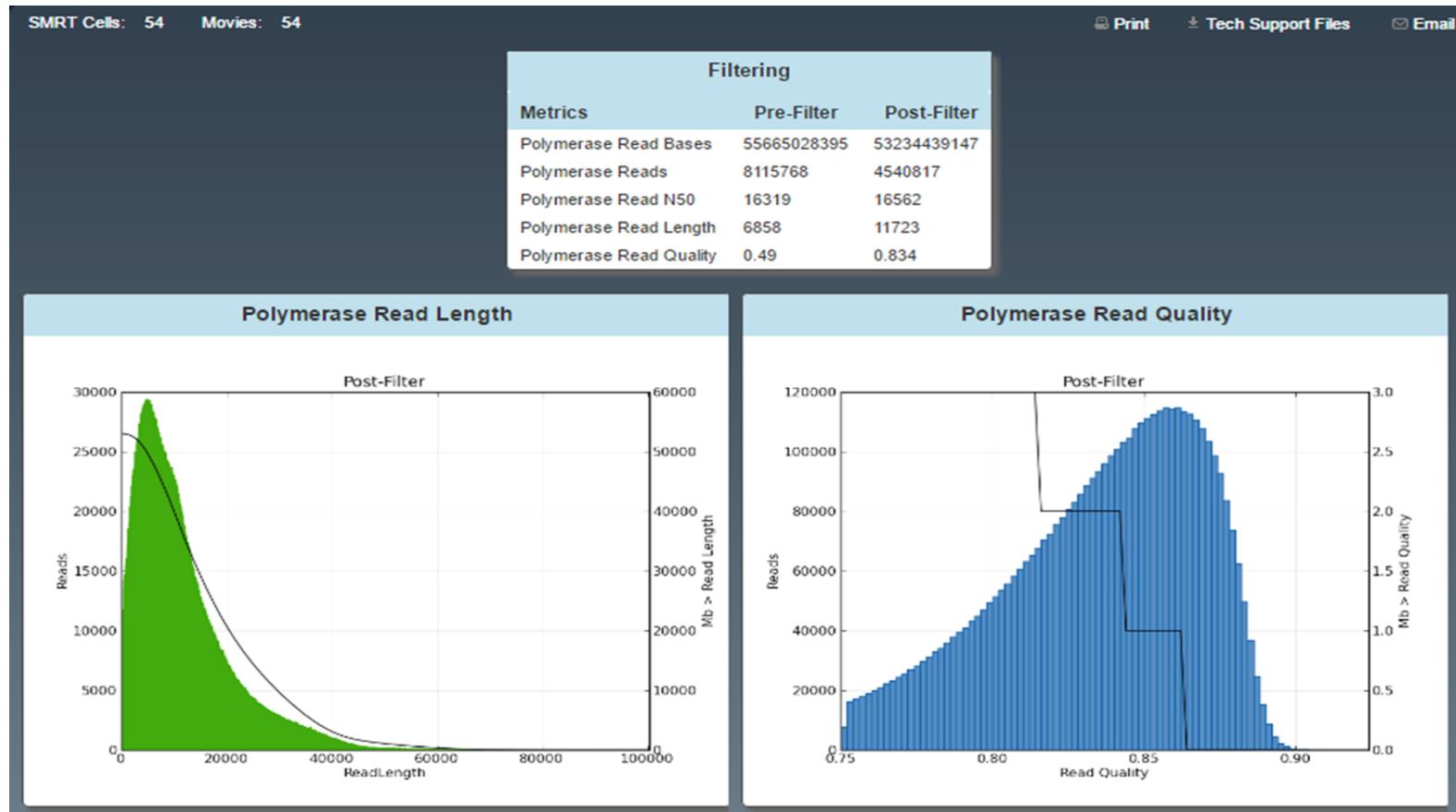


NC STATE UNIVERSITY
**Plants for
Human Health**
INSTITUTE

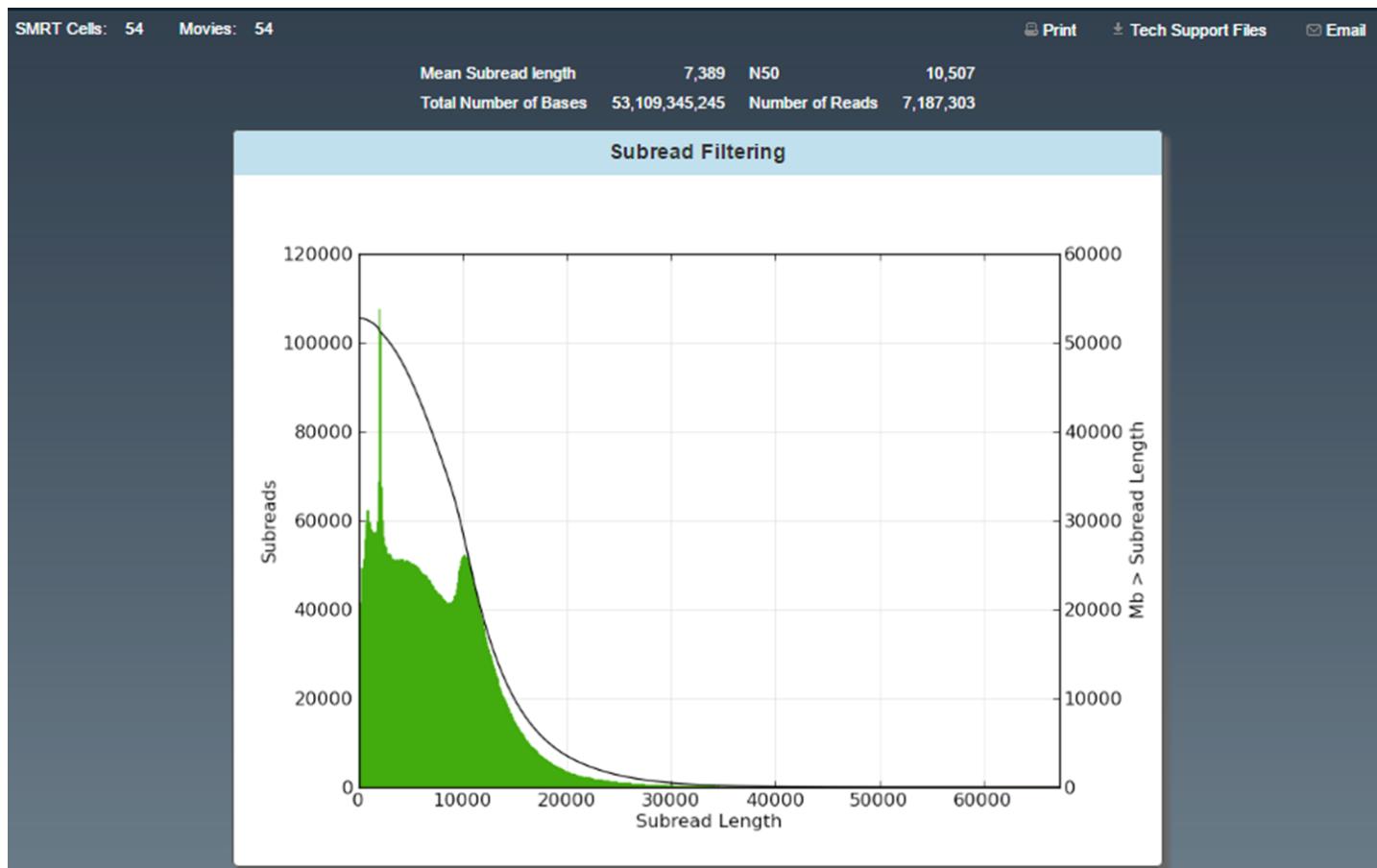
Blueberry Ploidy Level and Genome Size

- Naturally occurring 2X, 4X and 6X genomes
 - Most commercial varieties are 4X and 6X
 - Diploid blueberries are either ornamental or wild with no fresh or processed fruit commercial production use, but they have been used in breeding
- It is possible to cross a 2X with 4X and yet obtain a 4X progeny (unreduced gamete)
- It is possible to obtain a 5X genome by crossing a 4X and a 6X genome
 - But a 3X blueberry is rare or have not been successful
- Most breeding efforts are focused on 4X and 6X genomes
- It is believed that 4X blueberry is autotetraploid and 6X is a natural allohexaploid (2X + 4X ?)
- DNA content of a diploid genome is 2C ~1.37 pg, and the monoploid genome size is estimated ~670 Mb

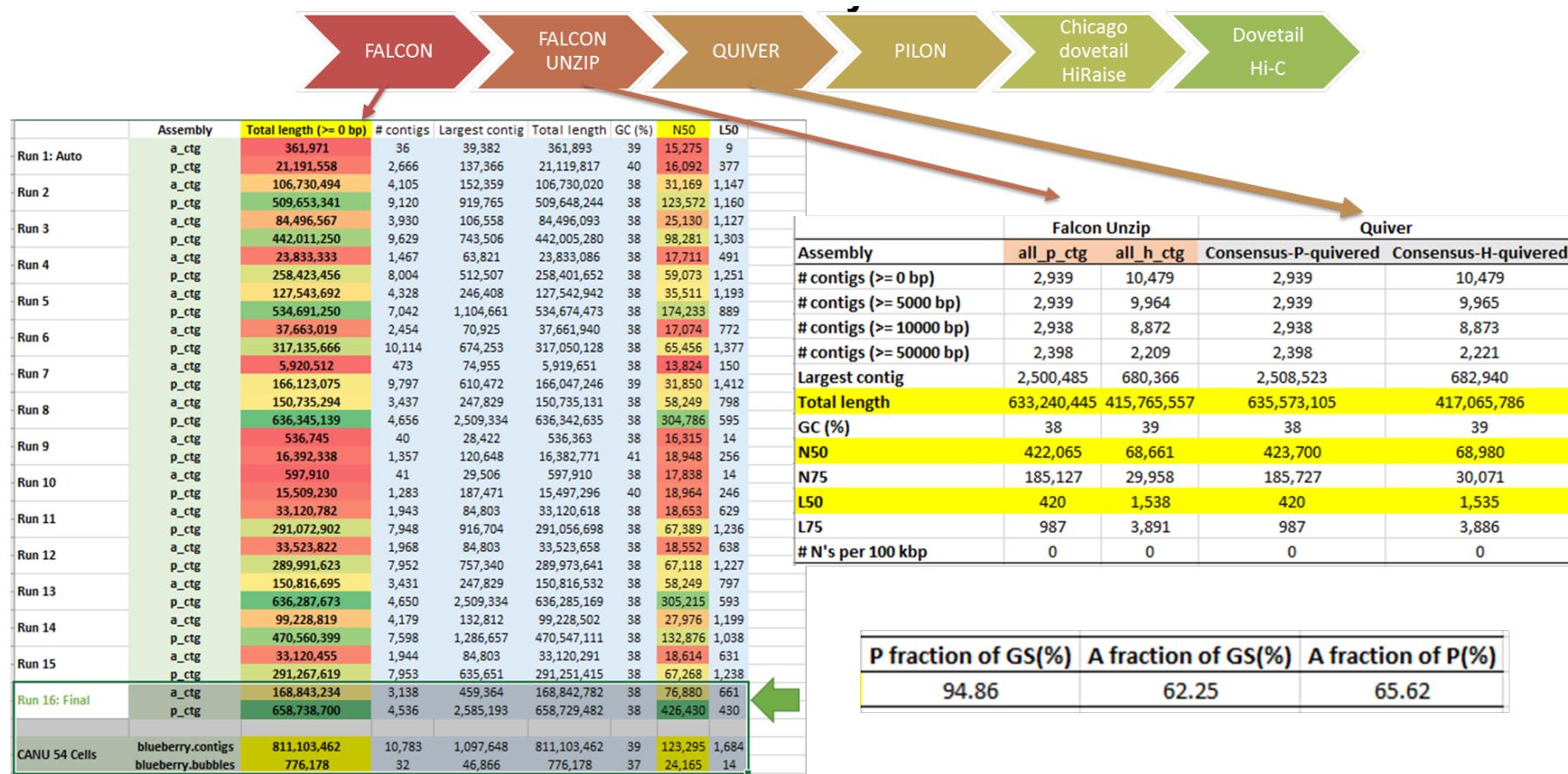
A Total of 54 RS II SMRT Cells Were Used for Blueberry Diploid Genome Assembly



A Total of 54 RS II SMRT Cells Were Used for Blueberry Diploid Genome Assembly



Blueberry Genome Assembly by Falcon – Falcon Unzip



code 3 (false insert on the genome and gap on the transcript)

000160F	647243	647244	000160F	647243	647264 000160F_pilon	647342 GAGAGAGAGAGAGAGAGAGA .	000160F	647173	647533
000160F	genome1	exon		647173	647533 93	+	.		
ID=c1138/f12p53/4042.mrna1.exon1;Name=c1138/f12p53/4042;Parent=c1138/f12p53/4042.mrna1;Target=c1138/f12p53/4042 1 339 + 1									

Note: after removing of the sequence on the unzip genome, the gap disappeared so the sequence was repeated 2 times.

False INSERTION causing gap

000160F
Sequence ID: Query_186133 Length: 764783 Number of Matches: 9
Range 1: 647173 to 656596

Score	Expect	Identities	Gaps	Strand	Frame
17219 bits(9324)	0.0()	9402/9433(99%)	31/9433(0%)	Plus/Plus	
Features:					
Query 1	CACTAGTCAAATGACAGCCGTTTTGATGCAACCATCGTTAGCCTCTCAGTATCTGAGT		60		
Sbjct 647173	CACTAGTCAAATGACAGCCGTTTTGATGCAACCATCGTTAGCCTCTCAGTATCTGAGT		647232		
Query 61	AATACTGTTT-----gagagagagagagagagagagagaga		98		
Sbjct 647233	AATACTGTTTGA		647292		
Query 99	gTAGTTGGAAAGGAAGTGAAGGAAAAAGGGATGGGATCCAAAAGAAGAGCTTGGTC		158		
Sbjct 647293	GTAGTTGGAAAGGAAGTGAAGGAAAAAGGGATGGGATCCAAAAGAAGAGCTTGGTC		647352		
Query 159	ATTCGATCCATTTCATGCATTGGCATGGTGAAGCTGCTAATGGGACTCGGATT		218		
Sbjct 647353	ATTCGATCCATTTCATGCATTGGCATGGTGAAGCTGCTAATGGGACTCGGATT		647412		
Query 219	CCTCGGGTCACTCGGCACGGGCTTCTATGCCGTATGTTGCTCTCAGCAAGAT		278		
Sbjct 647413	CCTCGGGTCACTCGGCACGGGCTTCTATGCCGTATGTTGCTCTCAGCAAGAT		647472		
Query 279	TATGAACATAATCGGGGGCGCTTCCGGGGGCCAAAACCTCTCATACCATCAACAA		338		
Sbjct 647473	TATGAACATAATCGGGGGCGCTTCCGGGGGCCAAAACCTCTCATACCATCAACAA		647532		
Query 339	GGTCCTCCCTCTCATAACTatataatataatataatataatataatataatata		398		
Sbjct 647533	GGTCCTCCCTCTCATAACTATATATATATATATATATATATATATATATATA		647592		
Query 399	tatataatatataatatataatatatacatcgatcatatgtatataatataCA		458		
Sbjct 647593	TATATATATATATATATATATATATATACACATCGTACATATGTATATATCA		647652		

000160F_pilon
Sequence ID: Query_179209 Length: 764876 Number of Matches: 9
Range 1: 647272 to 656682

Score	Expect	Identities	Gaps	Strand	Frame
17379 bits(9411)	0.0()	9411/9411(100%)	0/9411(0%)	Plus/Plus	
Features:					
Query 1	CACTAGTCAAATGACAGCCGTTTTGATGCAACCATCGTTAGCCTCTCAGTATCTGAGT		60		
Sbjct 647272	CACTAGTCAAATGACAGCCGTTTTGATGCAACCATCGTTAGCCTCTCAGTATCTGAGT		647331		
Query 61	AATACTGTTTgag		120		
Sbjct 647332	AATACTGTTTGA		647391		
Query 121	GAAAAGGGATGGGATCCAAAAGAAGAGCTTGGGTCAATTGATCCATTTCATGCAT		180		
Sbjct 647392	GAAAAGGGATGGGATCCAAAAGAAGAGCTTGGGTCAATTGATCCATTTCATGCAT		647451		
Query 181	TGGGATGGCATTGATAAGCTGCTAATGGGACTCGGATTCTCGGGTCACTCGGCAGCGG		240		
Sbjct 647452	TGGGATGGCATTGATAAGCTGCTAATGGGACTCGGATTCTCGGGTCACTCGGCAGCGG		647511		
Query 241	CTTCTATGCCGTATGTTGCTCTTCACTAGCAAGATATTGAACAATATCGGGGGCGCT		300		
Sbjct 647512	CTTCTATGCCGTATGTTGCTCTTCACTAGCAAGATATTGAACAATATCGGGGGCGCT		647571		
Query 301	TCCGCCGGGGCCAAAACCTCTCATACCATCAACAAAGGTCCCTCCCTCTCATAACTa		360		
Sbjct 647572	TCCGCCGGGGCCAAAACCTCTCATACCATCAACAAAGGTCCCTCCCTCTCATAACTa		647631		
Query 361	tatataatatataatatataatatataatatataatatataatatataatatata		420		
Sbjct 647632	TA		647691		
Query 421	tatataatatcacacatcgatcatatgtatataatCACACCGCAATGAACCTTGTACTA		480		
Sbjct 647692	TA		647751		

code 5 (false insert on the genome and stopping the transcript mapping to continue- partial mapping)

000886F	22018	22019	000886F	22018	22053	000886F_pilon	22009	CAACAACAACAAAACATAACCATAGTCCAAAGGGT	.	000886F	21276
	24514	000886F	genome1	exon	21276	24514	98	-	.		
ID=c18093/f2p28/4207.mrna1.exon1;Name=c18093/f2p28/4207;Parent=c18093/f2p28/4207.mrna1;Target=c18093/f2p28/4207 1 3258 + 1											

Note: in the unzip genome, the mapping stopped when reached to the sequence at base 22017 but in the pilon version the mapping continues since that sequence (interrupt) was removed so the gene is predicted complete.

000886F
Sequence ID: Query_210919 Length: 200276 Number of Matches: 4
Range 1: 15008 to 22017

Score	Expect	Identities	Gaps	Strand	Frame
12936 bits(7005)	0.0()	7010/7012(99%)	2/7012(0%)	Plus/Minus	

Features:

Query 2462	TGTTCAGTCCACTATGATGTCCTCATGAATCTCTTATCAGCTCGAATCTCAAAGAAC	2521
Sbjct 22017	TGTTCAGTCCACTATGATGTCCTCATGAATCTCTTATCAGCTCGAATCTCAAAGAAC	21958
Query 2522	AGTTGTTCAAGAAATCTATGGCATTCCACTTGGGTATAGCGGTTGTTGTC	2581
Sbjct 21957	AGTTGTTCAAGAAATCTATGGCATTTCGACTTGGGTATAGCGGTTGTTGTC	21898
Query 2582	TTGTTTCACTTCTCTCTAAAGACACTCACCCAGGCTCTGTTCTCATTATCAAATGCT	2641
Sbjct 21897	TTGTTTCACTTCTCTCTAAAGACACTCACCCAGGCTCTGTTCTCATTATCAAATGCT	21838
Query 2642	TTTGATATGTTCTTAGCCGATAAGTGAATTACTAACATTACCTGTGTTAAGTTCT	2701
Sbjct 21837	TTTGATATGTTCTTAGCCGATAAGTGAATTACTAACATTACCTGTGTTAAGTTCT	21778
Query 2702	ACCTTGTCTGACGGATGTGCACAATTATCCAACACTAGTATTAGTCTGGTCCACAG	2761
Sbjct 21777	ACCTTGTCTGACGGATGTGCACAATTATCCAACACTAGTATTAGTCTGGTCCACAG	21718
Query 2762	AACTCATGATGAGGCTTCTGTTCTTCAAGTATTGTGAAGAACGTGAAGTTGACTT	2821
Sbjct 21717	AACTCATGATGAGGCTTCTGTTCTTCAAGTATTGTGAAGAACGTGAAGTTGACTT	21658
Query 2822	CTTCAATAACTGCACTGATTGAGAAGACCATATCGGAGTGGTCTCTGTACCCCTCA	2881
Sbjct 21657	CTTCAATAACTGCACTGATTGAGAAGACCATATCGGAGTGGTCTCTGTACCCCTCA	21598
Query 2882	ATTTGAGGATAAGAGGACAAACATTGTAATTATCATTGATGAATGCAGCTTATAA	2941
Sbjct 21597	ATTTGAGGATAAGAGGACAAACATTGTAATTATCATTGATGAATGCAGCTTATAA	21538
Query 2942	CTTAAGTCCCTTCACATTGGTGGAAAGAGATGTTCAATCAAAGATTGGAAAATGTGG	3001
Sbjct 21537	CTTAAGTCCCTTCACATTGGTGGAAAGAGATGTTCAATCAAAGATTGGAAAATGTGG	21478

000886F_pilon

Sequence ID: Query_245069 Length: 200282 Number of Matches: 3
Range 1: 14997 to 24469

Score	Expect	Identities	Gaps	Strand	Frame
17494 bits(9473)	0.0()	9473/9473(100%)	0/9473(0%)	Plus/Minus	

Features:

Query 1	CTTGTCAAAACGGGGCGGTGCGATCTCCAACTGGTATAACCGTCGCATACGAGACGAA	60
Sbjct 24469	CTTGTCAAAACGGGGCGGTGCGATCTCCAACTGGTATAACCGTCGCATACGAGACGAA	24410
Query 61	TCTGTAAAACAGATAAGGATCTGAGGTTGAAATCTTGTGAAATTGACATTCACGATCAG	120
Sbjct 24409	TCTGTAAAACAGATAAGGATCTGAGGTTGAAATCTTGTGAAATTGACATTCACGATCAG	24350
Query 121	GTAAAACGACAAACGCTGAATTAGCATATTGTGTCGATTATCACGCTTAAAGGTCATT	180
Sbjct 24349	GTAAAACGACAAACGCTGAATTAGCATATTGTGTCGATTATCACGCTTAAAGGTCATT	24290
Query 181	ATTTAGGTTGCTGTAATTAGATGATTGTCGTGTCGACTAGTCATCGTGTGATTAA	240
Sbjct 24289	ATTTAGGTTGCTGTAATTAGATGATTGTCGTGTCGACTAGTCATCGTGTGATTAA	24230
Query 241	GCGCATATGTTGTCATTACAAGTTGTGAACTAGGATATTGAGGTGCTAGGGTTGAGT	300
Sbjct 24229	GCGCATATGTTGTCATTACAAGTTGTGAACTAGGATATTGAGGTGCTAGGGTTGAGT	24170
Query 301	TGATTCTGAAAGCATTAGATTGGCCGCTTTGATTGAGCCTAGGTGAGTTGGCG	360
Sbjct 24169	TGATTCTGAAAGCATTAGATTGGCCGCTTTGATTGAGCCTAGGTGAGTTGGCG	24110
Query 361	TTTCTTAGTTGAGGTTCAATTGATGCCAAGTATTAGTTGGGTTGGGCT	420
Sbjct 24109	TTTCTTAGTTGAGGTTCAATTGATGCCAAGTATTAGTTGGGTTGGGCT	24050
Query 421	AATTACTGCTGAGGAAACAATCGTGGTGAATTCTATGTCCTTGGTGGCAG	480
Sbjct 24049	AATTACTGCTGAGGAAACAATCGTGGTGAATTCTATGTCCTTGGTGGCAG	23990
Query 481	GGGAGGATTGGGAATTGAGTTGGGTGTGTTCTGATGGGGCACTGAAACTGAGTTGG	540

False INSERTION causing stop

code 6 (false deletion on the genome and gap on the genome when mapping the transcript)

000104F	111391	111392	000104F	111391 000104F_pilon	111417	111427 . TCTCTCTCT	000104F	110997	111584	000104F
genome1	exon		110997	111584	98	.				
ID=c13750/f25p70/3400.mrna1.exon1;Name=c13750/f25p70/3400;Parent=c13750/f25p70/3400.mrna1;Target=c13750/f25p70/3400 1 591 +										1

Note: on the unizp genome, this sequence is missing which caused a gap on position 111391, when the sequence is inserted/corrected on the pilon genome, the gap is removed and replaced with the matches.

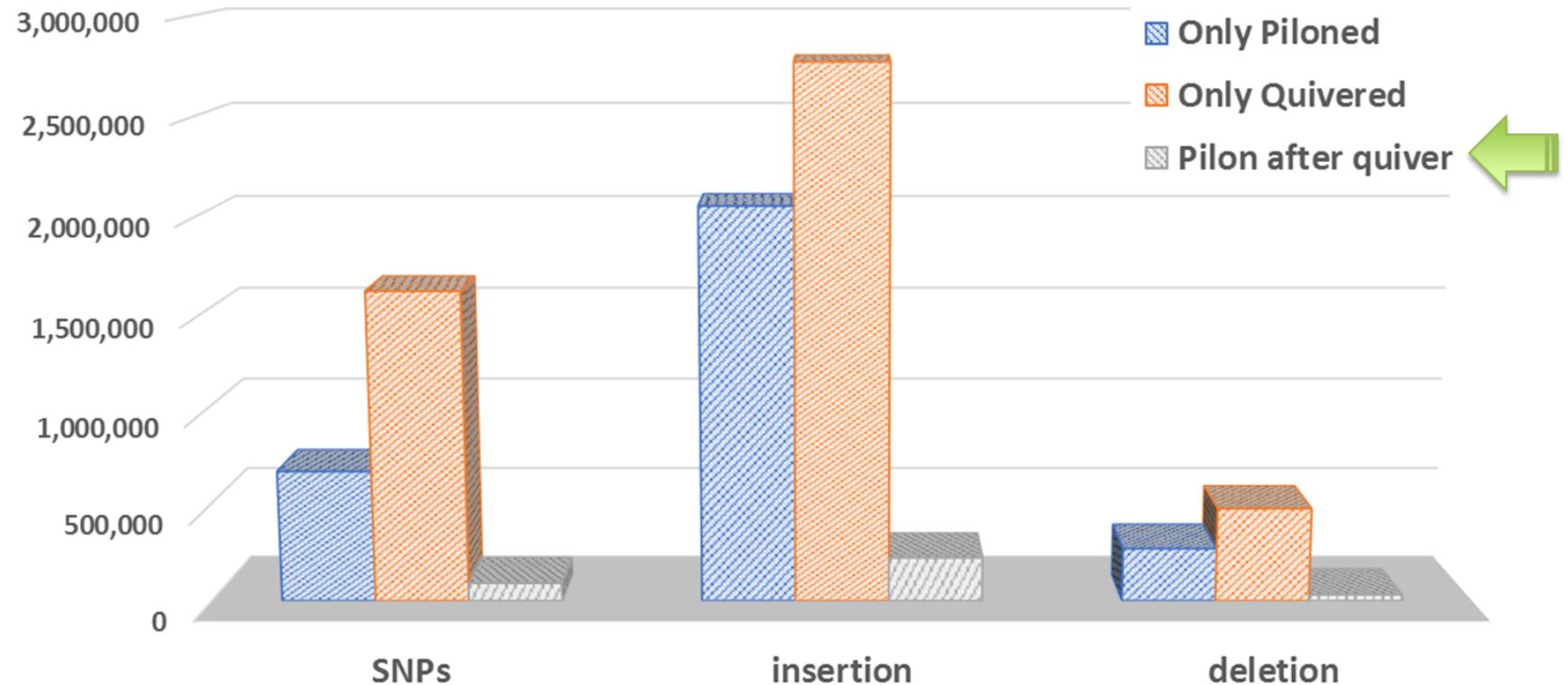
000104F
Sequence ID: Query_213451 Length: 924261 Number of Matches: 5
Range 1: 99245 to 111584

Score	Expect	Identities	Gaps	Strand	Frame
22718 bits(12302)	0.0()	12340/12355(99%)	15/12355(0%)	Plus/Minus	
Features:					
Query 1	AAGTTGAAAAATACAGCTCATTGTGTTTCAAGAAAAGAGGTTGCAAATCCAGCTATTAC	60			
Sbjct 111584	AAGTTGAAAAATACAGCTCATTGTGTTTCAAGAAAAGAGGTTGCAAATCCAGCTATTAC	111525			
Query 61	AGGACACTGGAAAAAAATCCTCACTGTGATCATTCACAAATCCAAATCCAAATTGCC	120			
Sbjct 111524	AGGACACTGGAAAAAAATCCTCACTGTGATCATTCACAAATCCAAATCCAAATTGCC	111466			
Query 121	ATAAGATCTTCCCCCTCGCCGCCCTAATCTTGACCTTTGCAATCTATGTagagaga	180			
Sbjct 111465	ATAAGATCTTCCCCCTCGCCGCCCTAATCTTGACCTTTGCAATCTATGTagAGAG-	111407			
Query 181	gagagagagagagagagagagagagagGGTCCGTTTAGTTGGGAGAGCCCACC	240			
Sbjct 111406	-----GAG-GAGAGAGAGAG-GAGAGGGTCCGTTTAGTTGGGAGAGCCCACC	111359			
Query 241	ACCAAGAGACGTGGTATTGCTGATCTACCTTTACTTATCTTGATCAGATAGATT	300			
Sbjct 111358	ACCAAGAGACGTGGTATTGCTGATCTACCTTTACTTATCTTGATCAGATAGATT	111299			
Query 301	CTTTTGGTTCTCTGTAGCCATTGACTCTTGGACTGGGAATTGAGAAAATTGGGATA	360			
Sbjct 111298	CTTTTGGTTCTCTGTAGCCATTGACTCTTGGACTGGGAATTGAGAAAATTGGGATA	111239			
Query 361	TTTGTGTTGGTGGTGGTTGGCTCAGTCTAGATGGGCTGGAAGCCGTTTCTA	420			
Sbjct 111238	TTTGTGTTGGTGGTGGTTGGCTCAGTCTAGATGGGCTGGAAGCCGTTTCTA	111179			
Query 421	TCCGGCTCCGACCTATCGTCTTTAGAGACTTACTGGGACACCGACGACGCCCTGG	480			
Sbjct 111178	TCCGGCTCCGACCTATCGTCTTTAGAGACTTACTGGGACACCGACGACGCCCTGG	111119			
Query 481	CCCACGGCTGGGCCACACTCTCACCGCCATCGCCGCTACTAAACCCACGGCCCCCGCT	540			
Sbjct 111144	CCCACGGCTGGGCCACACTCTCACCGCCATCGCCGCTACTAAACCCACGGCCCCCGCT	111025			

000104F_pilon
Sequence ID: Query_221427 Length: 924449 Number of Matches: 5
Range 1: 99270 to 111624

Score	Expect	Identities	Gaps	Strand	Frame
22816 bits(12355)	0.0()	12355/12355(100%)	0/12355(0%)	Plus/Minus	
Features:					
Query 1	AAGTTGAAAAATACAGCTCATTGTGTTTCAAGAAAAGAGGTTGCAAATCCAGCTATTAC	60			
Sbjct 111624	AAGTTGAAAAATACAGCTCATTGTGTTTCAAGAAAAGAGGTTGCAAATCCAGCTATTAC	111565			
Query 61	AGGACACTGGAAAAAAATCCTCACTGTGATCATTCACAAATCCAAATCCAAATTGCC	120			
Sbjct 111564	AGGACACTGGAAAAAAATCCTCACTGTGATCATTCACAAATCCAAATCCAAATTGCC	111505			
Query 121	ATAAGATCTTCCCCCTCGCCGCCCTAATCTTGACCTTTGCAATCTATGTagagaga	180			
Sbjct 111504	ATAAGATCTTCCCCCTCGCCGCCCTAATCTTGACCTTTGCAATCTATGTagAGAGAGA	111445			
Query 181	gagagagagagagagagagagagagagGGTCCGTTTAGTTGGGAGAGCCCACC	240			
Sbjct 111444	GAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGGGTCCGTTTAGTTGGGAGAGCCCACC	111385			
Query 241	ACCAAGAGACGTGGTATTGCTGATCTACCTTTACTTATCTTGATCAGATAGATT	300			
Sbjct 111384	ACCAAGAGACGTGGTATTGCTGATCTACCTTTACTTATCTTGATCAGATAGATT	111325			
Query 301	CTTTTGGTTCTCTGTAGCCATTGACTCTTGGACTGGGAATTGAGAAAATTGGGATA	360			
Sbjct 111324	CTTTTGGTTCTCTGTAGCCATTGACTCTTGGACTGGGAATTGAGAAAATTGGGATA	111265			
Query 361	TTTGTGTTGGTGGTGGTTGGCTCAGTCTAGATGGGCTGGAAGCCGTTTCTA	420			
Sbjct 111264	TTTGTGTTGGTGGTGGTTGGCTCAGTCTAGATGGGCTGGAAGCCGTTTCTA	111205			
Query 421	TCCGGCTCCGACCTATCGTCTTTAGAGACTTACTGGGACACCGACGACGCCCTGG	480			
Sbjct 111204	TCCGGCTCCGACCTATCGTCTTTAGAGACTTACTGGGACACCGACGACGCCCTGG	111145			
Query 481	CCCACGGCTGGGCCACACTCTCACCGCCATCGCCGCTACTAAACCCACGGCCCCCGCT	540			
Sbjct 111144	CCCACGGCTGGGCCACACTCTCACCGCCATCGCCGCTACTAAACCCACGGCCCCCGCT	111025			

Comparison of Base Correction Using Different Methods



Worth to mention that:

The total number of bases Pilon corrected was < 0.07% of the total genome size assembled

So if believe in Pilon, after Falcon assembly and Quivering, we already have a genome assembled of over 99.93% accuracy.

Iso-Seq Project

Things to consider

1. Choice germplasm source for RNA extraction
 - Diploid (for genome annotation as well as comparison)
 - Tetraploid
 - Hexaploid
2. Variety or cultivar
 - **W85-20** *V. caesariense* (N.J. blueberry)
 - **O'Neal** *V. corymbosum* (native to north east of the U.S.)
 - **Premier** *V. virginatum* or *ashei* (native to southern eastern U.S.)
3. Tissues

Iso-Seq Project

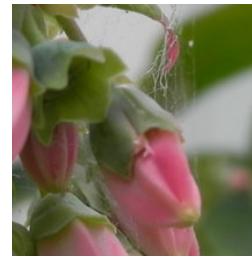
Tissue Samples Were Collected from Field Grown Blueberries; Leaf, Flower, Fruit, Root



Leaf



Flower Stage 1



Flower Stage 3



Flower Stage 5



Fruit Stage 1



Fruit Stage 2



Fruit Stage 3



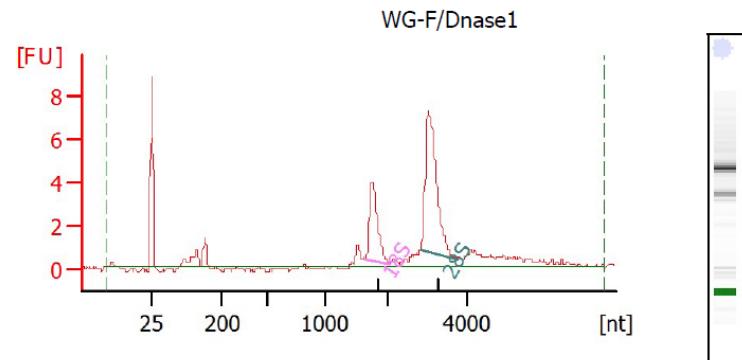
Fruit Stage 4



Root

Total RNA Extraction Was Attempted by Different Kits

- Sigma Plant RNA extraction Kit
- Bioanalyzer was used to check the quality of RNA
- The same RNA for both Illumina and Iso-Seq libraries
- Iso-Seq libraries were made with size selection option
- KAPA stranded RNA-Seq Kit to make Illumina libraries



Overall Results for sample 3 : WG-F/Dnase1

RNA Area: 46.8
RNA Concentration: 36 ng/ μ l
rRNA Ratio [28s / 18s]: 2.3
RNA Integrity Number (RIN): 9.7 (B.02.08)
Result Flagging Color: 
Result Flagging Label: RIN: 9.70

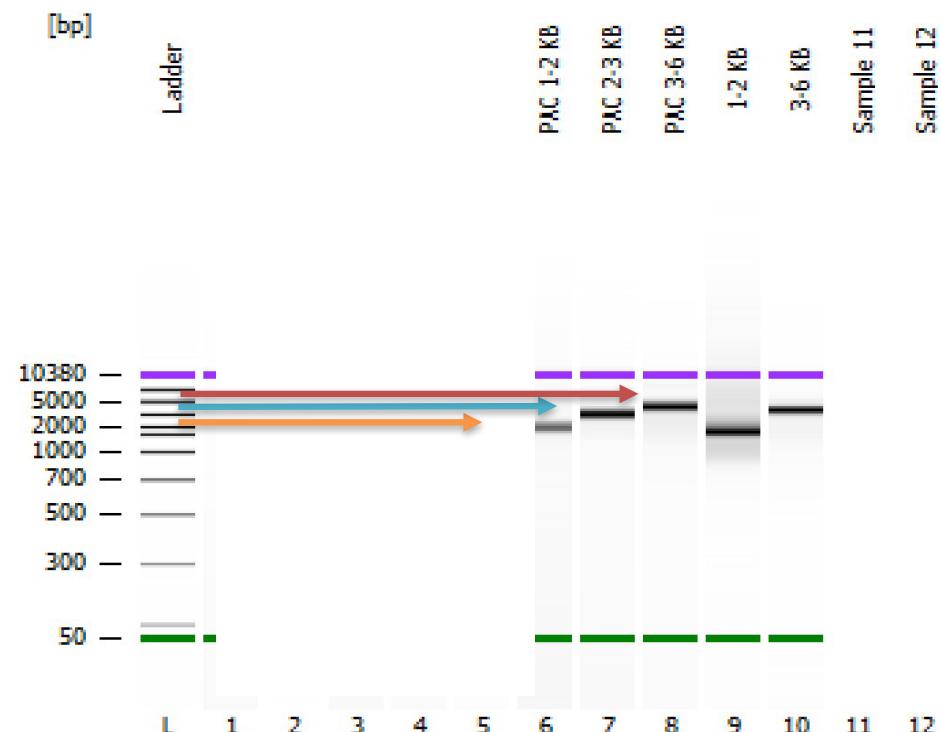
Fragment table for sample 3 : WG-F/Dnase1

Name	Start Size [nt]	End Size [nt]	Area	% of total Area
------	-----------------	---------------	------	-----------------

18S	1,651	2,091	7.2	15.4
28S	2,836	3,733	16.2	34.7

Barcoded Library Construction for PacBio Iso-seq Sequencing

BluePippin Size Selection



Making Barcoded Iso-Seq Libraries

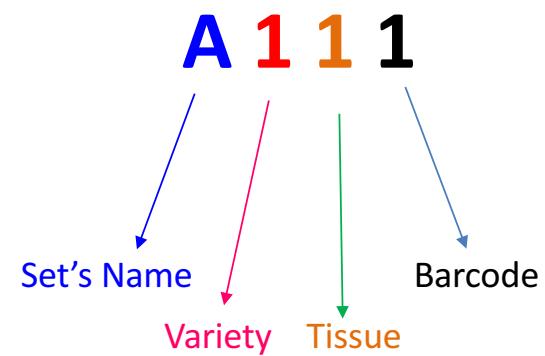
- There were (are) six barcoded adapters available to make pooled libraries

	Primer Sequence	16-mer barcode	oligo dT
dT_BC1	AAGCAGTGGTATCAACGCAGAGTAC	ttagacgtgcgtcatTTTTTTTTTTTTTTTTTTTTTTTTTTVN	
dT_BC2	AAGCAGTGGTATCAACGCAGAGTAC	ctatacatgactctgcTTTTTTTTTTTTTTTTTTTTTTTTVN	
dT_BC3	AAGCAGTGGTATCAACGCAGAGTAC	tactagagtagcactcTTTTTTTTTTTTTTTTTTTTVN	
dT_BC4	AAGCAGTGGTATCAACGCAGAGTAC	tgtgtatcgtacatgTTTTTTTTTTTTTTTTTTVN	
dT_BC5	AAGCAGTGGTATCAACGCAGAGTAC	gatctctactatatgcTTTTTTTTTTTTTTTTVN	
dT_BC1	AAGCAGTGGTATCAACGCAGAGTAC	acagtctatactgctgTTTTTTTTTTTTTTTTVN	

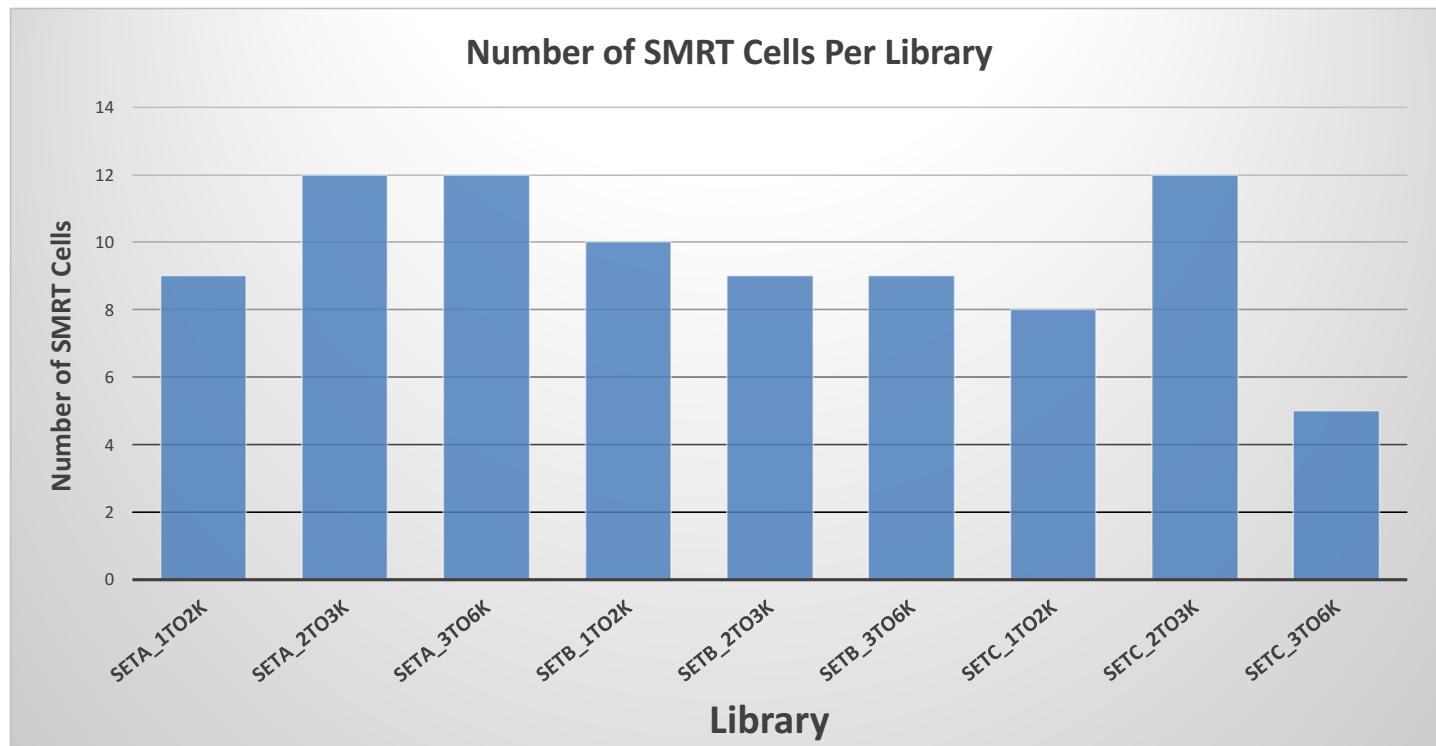
- Nine tissues for each variety (4X and 6X)
 - Only 4X and 6X libraries were barcoded (RSII Sequencing)
 - 2X library was not barcoded (2 Sequel SMRT cells to date)

We defined sets (A, B and C)

- Our approach
 - Make groups of six and having an indexing convention
 - A111, A122, A133, A144, A155, A166 (1-2, 2-3, 3-6 kb)
 - B171, B182, B193, B214, B225, B236 (1-2, 2-3, 3-6 kb)
 - C241, C252, C263, C274, C285, C296 (1-2, 2-3, 3-6 kb)
 - SMRT analysis 2.3.0 to analyze the data

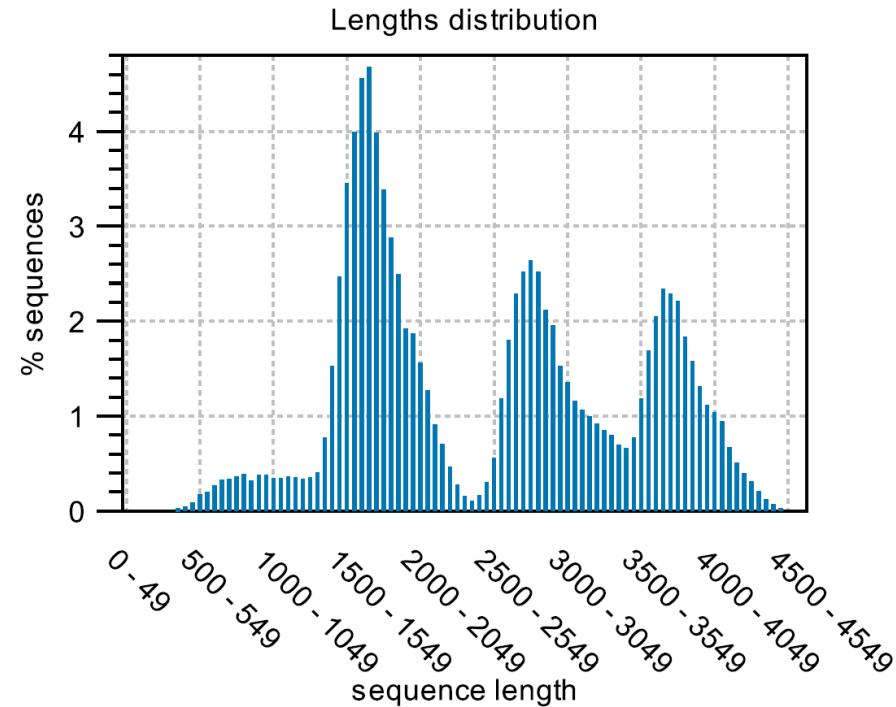


On Average We Ran 9.5 SMRT Cells Per Library (87 RSII total)



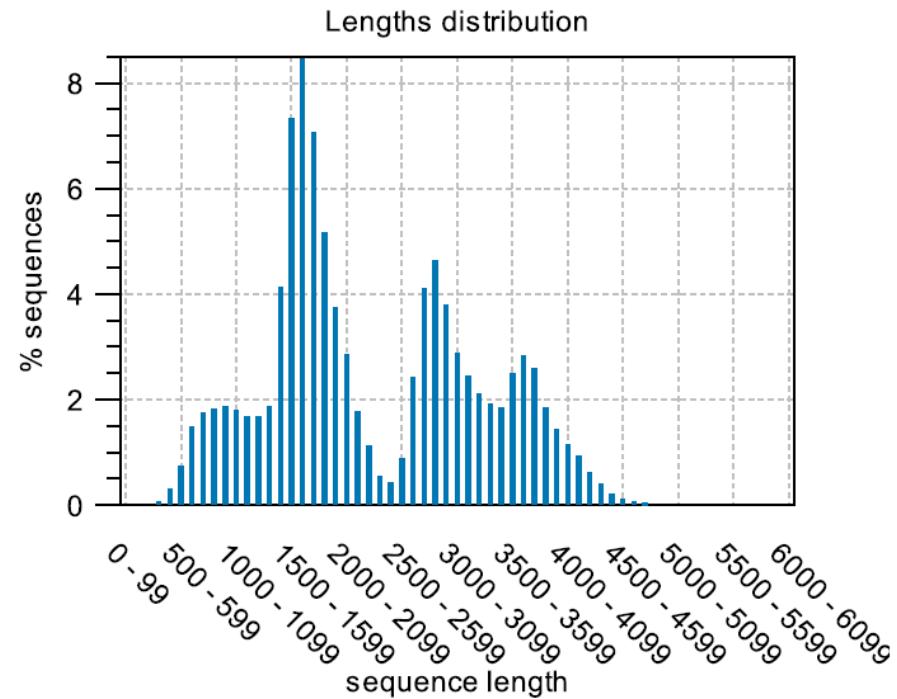
Quiver Was Used to Obtain HQ and LQ Data

- Full Length Non-Chimeric (FLNC) = 1,624,690
- Tetraploid Genome LQ = 773,571
- Tetraploid Genome HQ = 141,399 (351,343,616 nt)



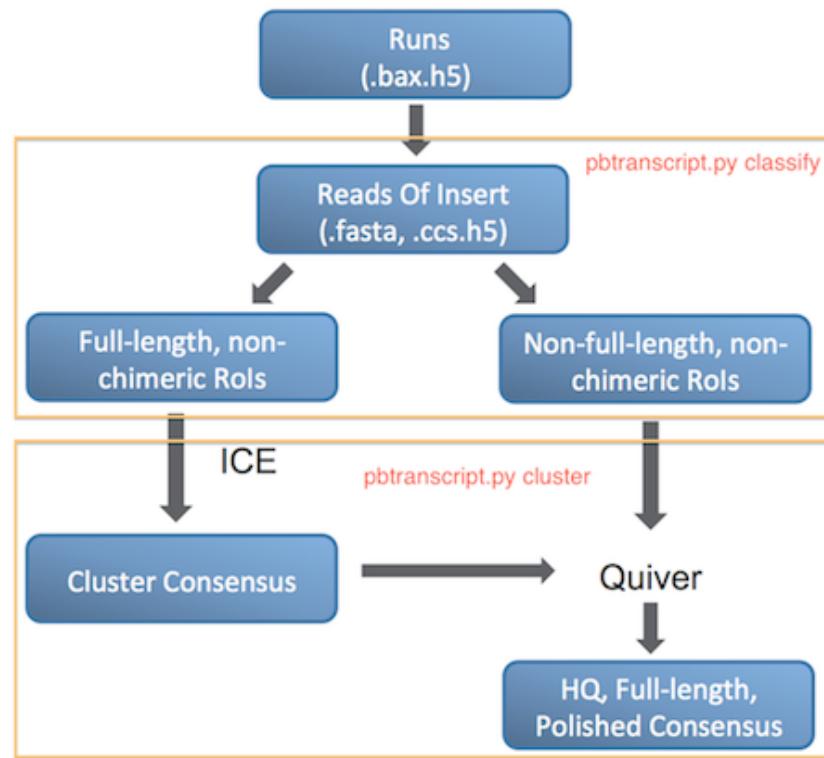
Quiver Was Used to Obtain HQ and LQ Data

- Full Length Non-Chimeric (FLNC) = 1,302,432
- Hexaploid Genome LQ = 614,378
- Hexaploid genome HQ = 110,050 and (250,964,213 nt)



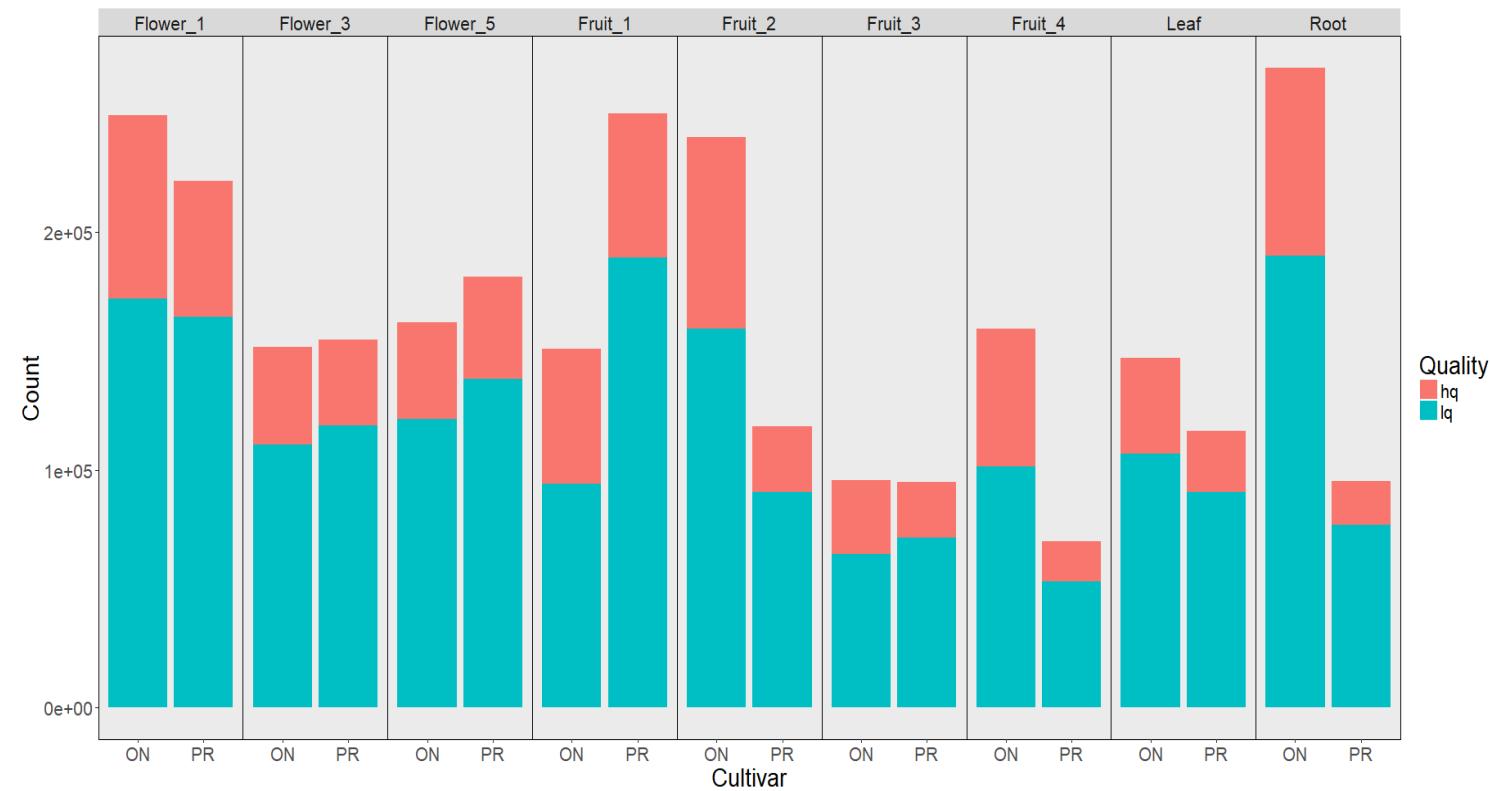
Data analysis pipeline (RSII Data)

All SMRT cells of each cultivar were used to output FLNC, LQ and HQ reads



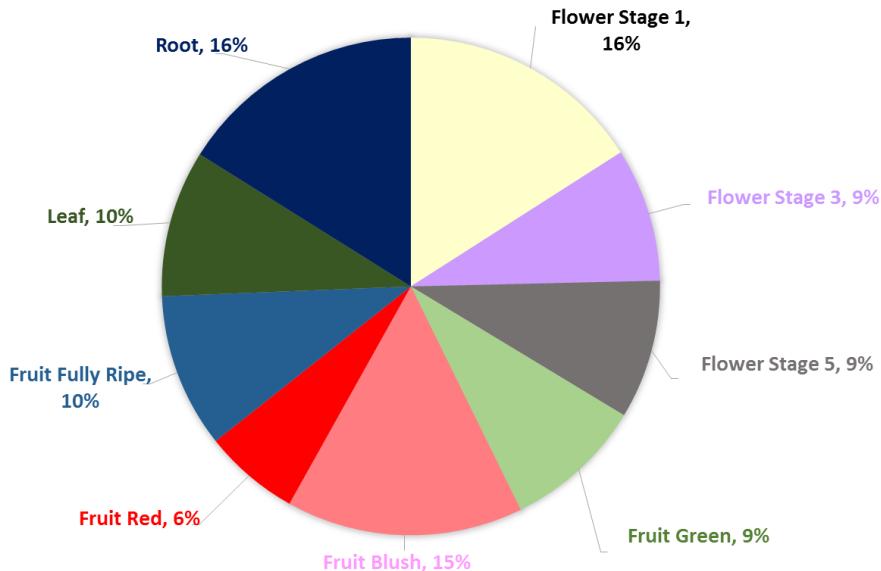
https://github.com/PacificBiosciences/cDNA_primer/wiki/RS_IsoSeq-%28v2.3%29-Tutorial-%231.-Getting-full-length-reads

Proportion of LQ and HQ Reads in 4X and 6X Blueberry Genotypes Separated by Different Tissues

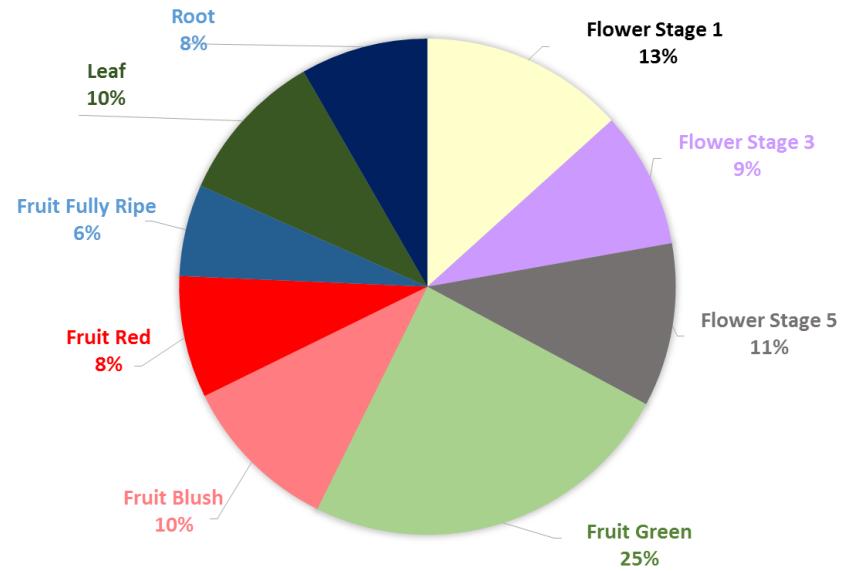


Demultiplexing Full Length Non-chimeric Reads (FLNC)

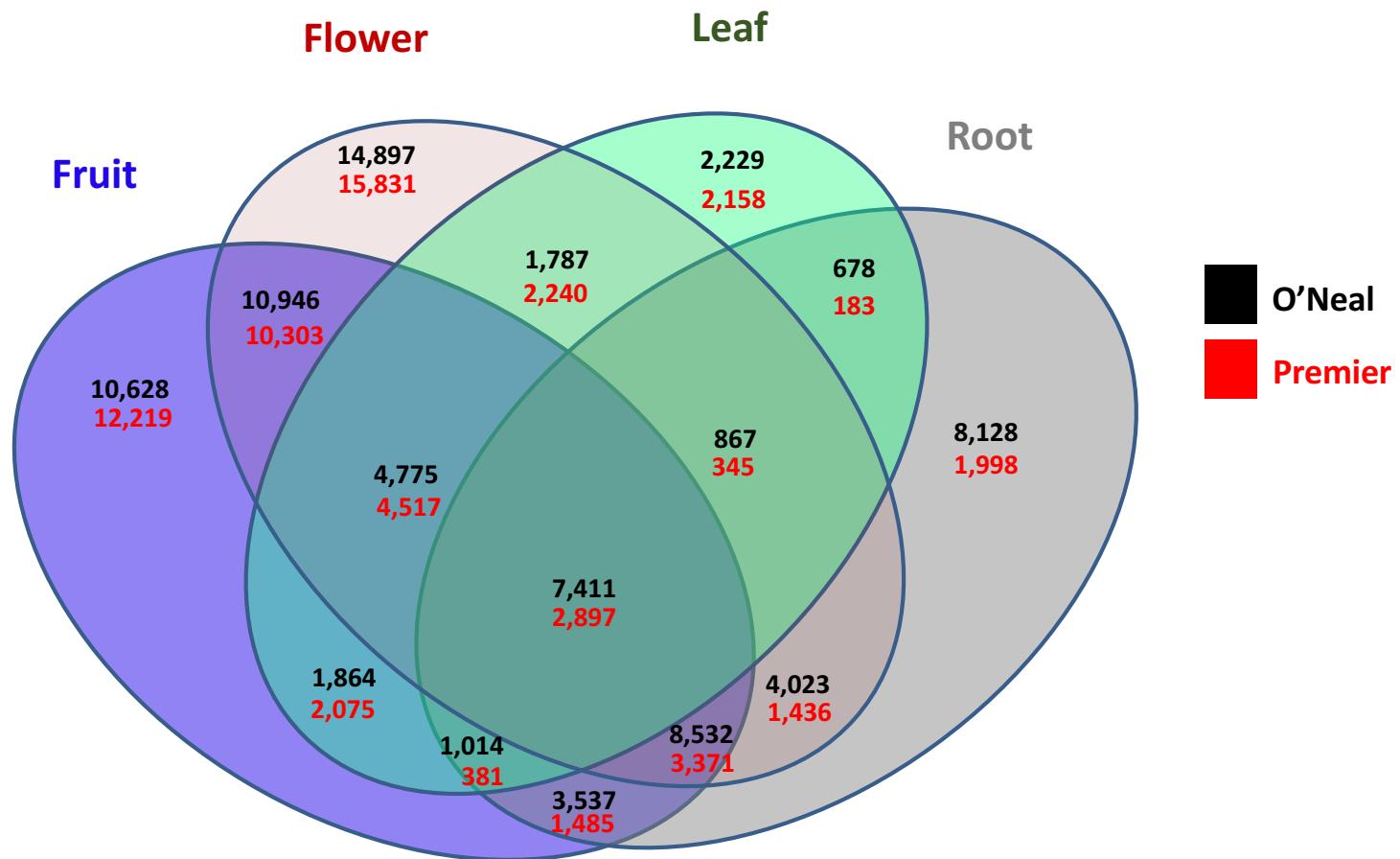
PERCENT OF BARCODES FOR EACH TISSUE TYPE OF O'NEAL



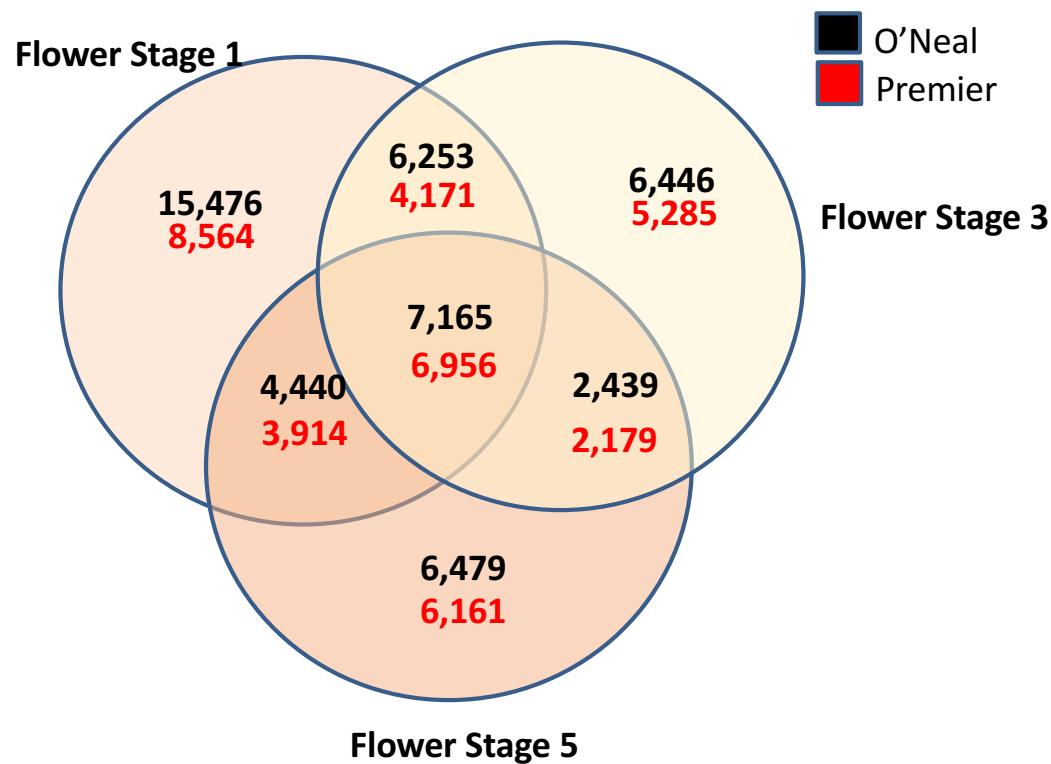
PERCENT OF BARCODES FOR EACH TISSUE TYPE OF PREMIER



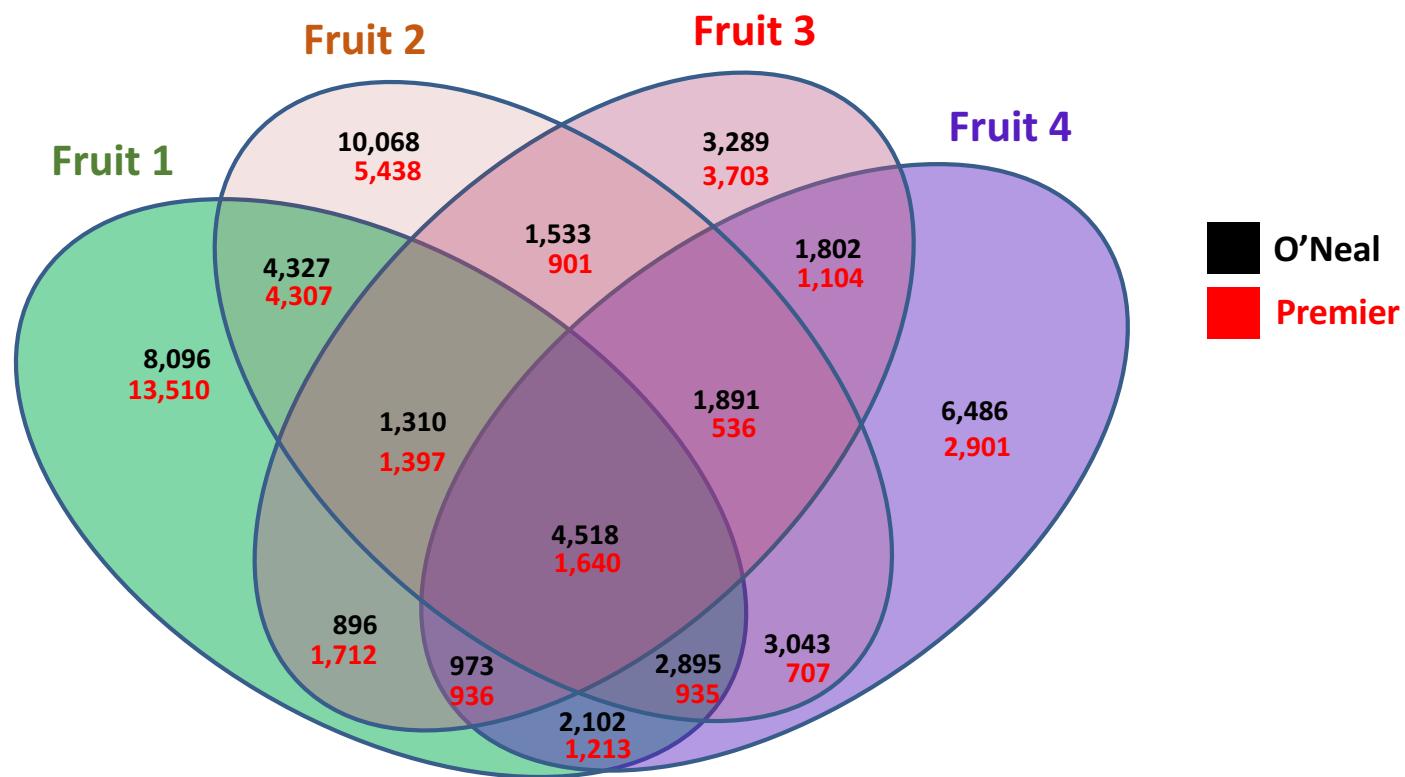
Overlap of all PacBio Isoforms in Four Tissues of 4X and 6X Genotypes



Overlap of all PacBio Isoforms in Three Flower Development Stages of 4X and 6X Genotypes



Number of Iso-Seq Sequences Separated by Tissue Type in 4X and 6X Genomes



GMAP Was Used to Map HQ Reads to the Reference Sequence

- Two versions of reference sequence was used
 - Quivered genome
 - Pilon corrected Quivered genome
(<https://github.com/broadinstitute/pilon/wiki>)
- Mapping was done to primary contigs
- It ran once with default parameters only for 4X and 6X iso-seq data
- It ran for the second time with the following parameters for 2x, 4X and 6X
 - $--min-trimmed-coverage = 0.95$ & $--min-identity=0.95$

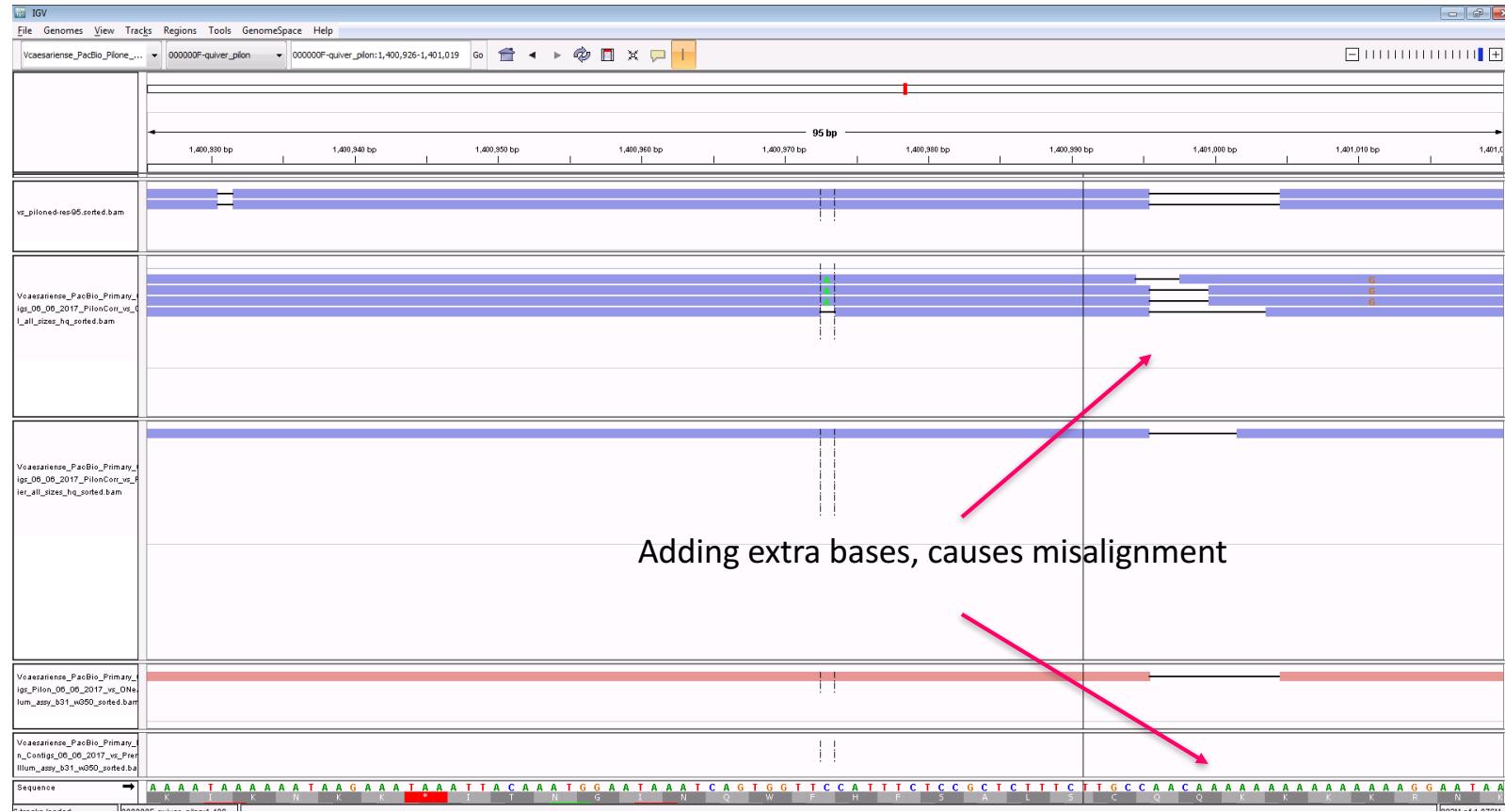
GMAP Default Values					
No. of Iso-Seqs Mapped to Quivered Genome			No. of Iso-Seqs Mapped to Quivered – Pilonized Genome		
2X	4X	6X	2X	4X	6X
-	140,664 (99.48%)	109,672 (99.63%)	-	140,096 (99.07%)	109,099 (99.13%)
				 568	 573
No. of Iso-Seqs Uniquely Mapped to Quivered Genome			No. of Iso-Seqs Mapped to Quivered – Pilonized Genome		
-	100,441	71,146	-	103,106  2,665	77,613  6,467

The increase in the number of uniquely mapped sequences to the [Pilon corrected genome](#), may indicate that it is better to make this correction.

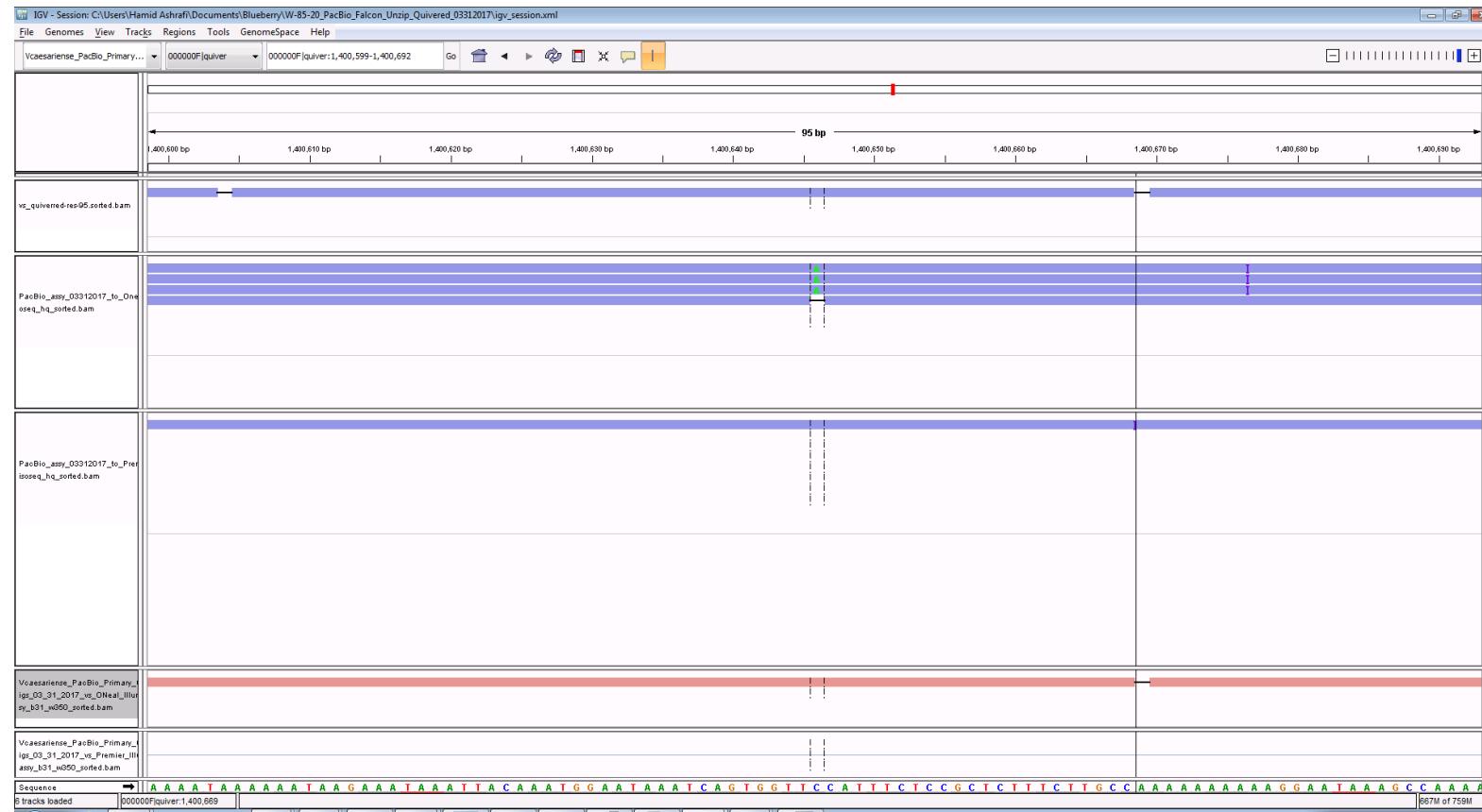
GMAP Run With min cov 0.95 and min identity 0.95

No. of Iso-Seqs Mapped to Quivered Genome			No. of Iso-Seqs Mapped to Quivered – Pilonized Genome		
2X	4X	6X	2X	4X	6X
-	127,312 (90.00%)	96,681 (87.85%)	-	127,307 5 (90.00%)	96,709 28 (87.87%)
No. of Iso-Seqs Uniquely Mapped to Quivered Genome			No. of Iso-Seqs Mapped to Quivered – Pilonized Genome		
29,512	100,428	74,097	29,509 3	100,329 99	74,074 23

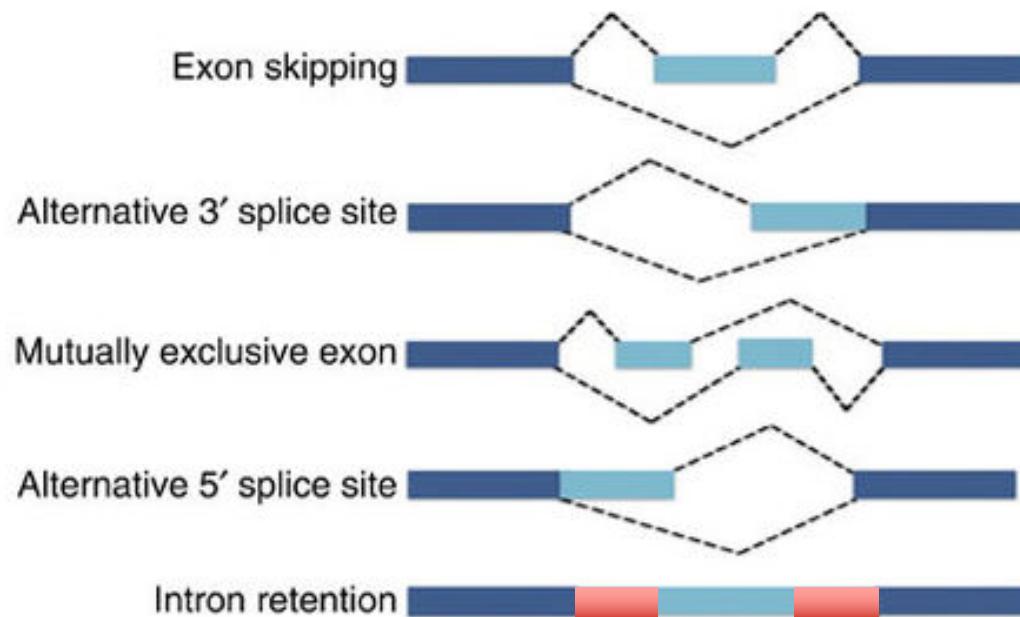
Should We Correct The Genome with Pilon?



Aligning Iso-Seq to Only Quivered Genome



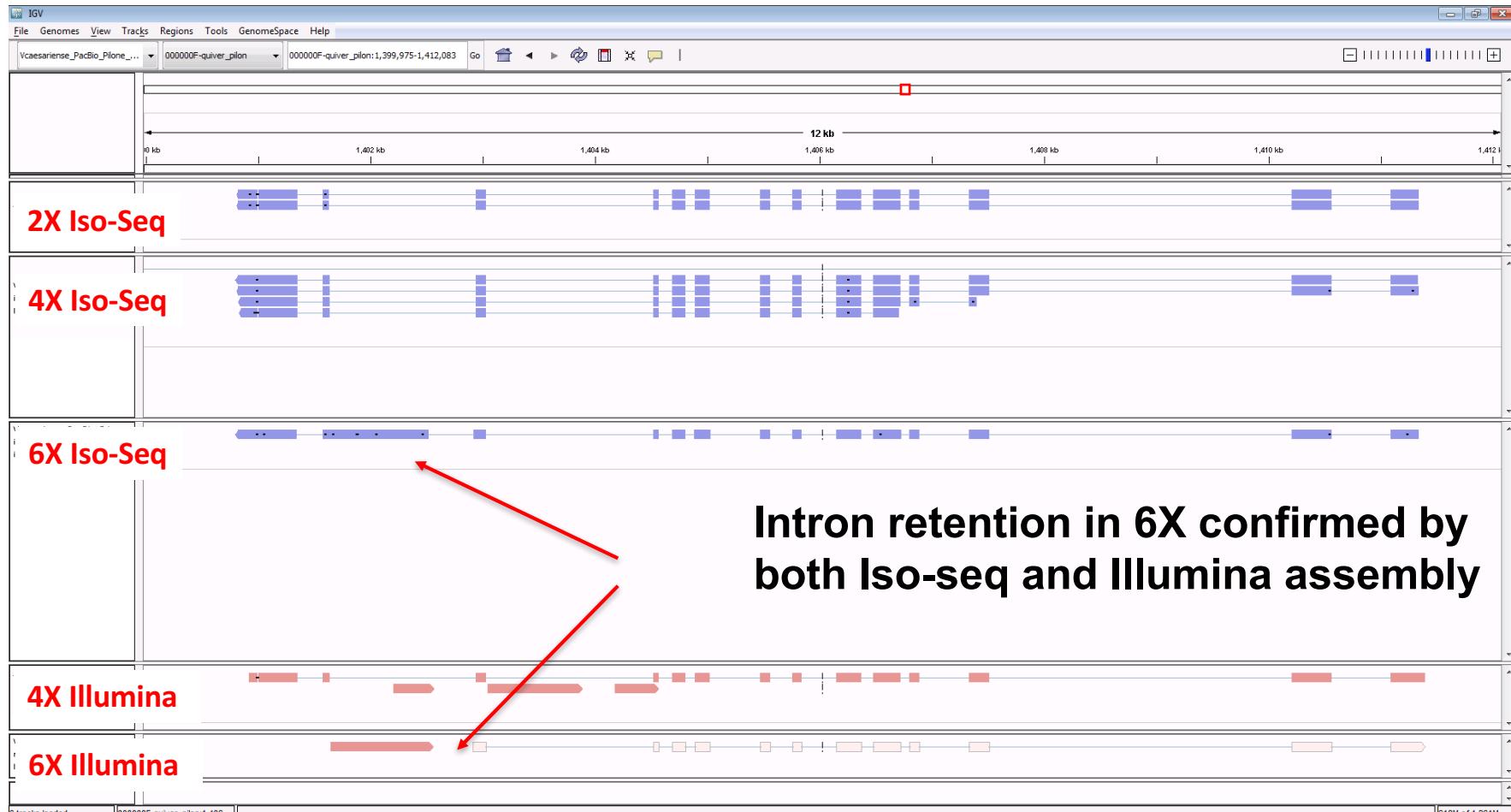
Five Alternative Splicing Modes

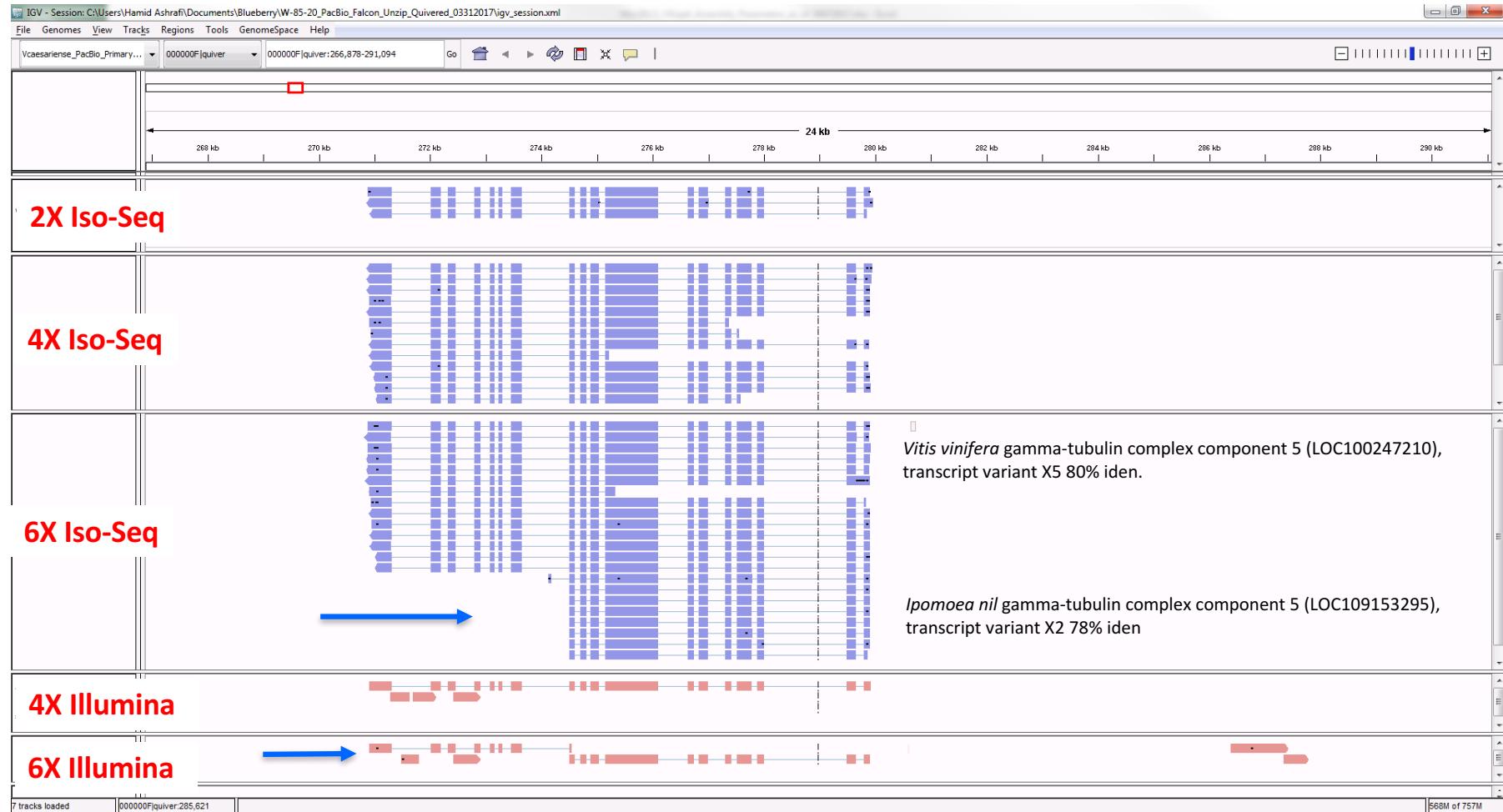


Wang et al. (2016) <https://www.nature.com/articles/ncomms11708>

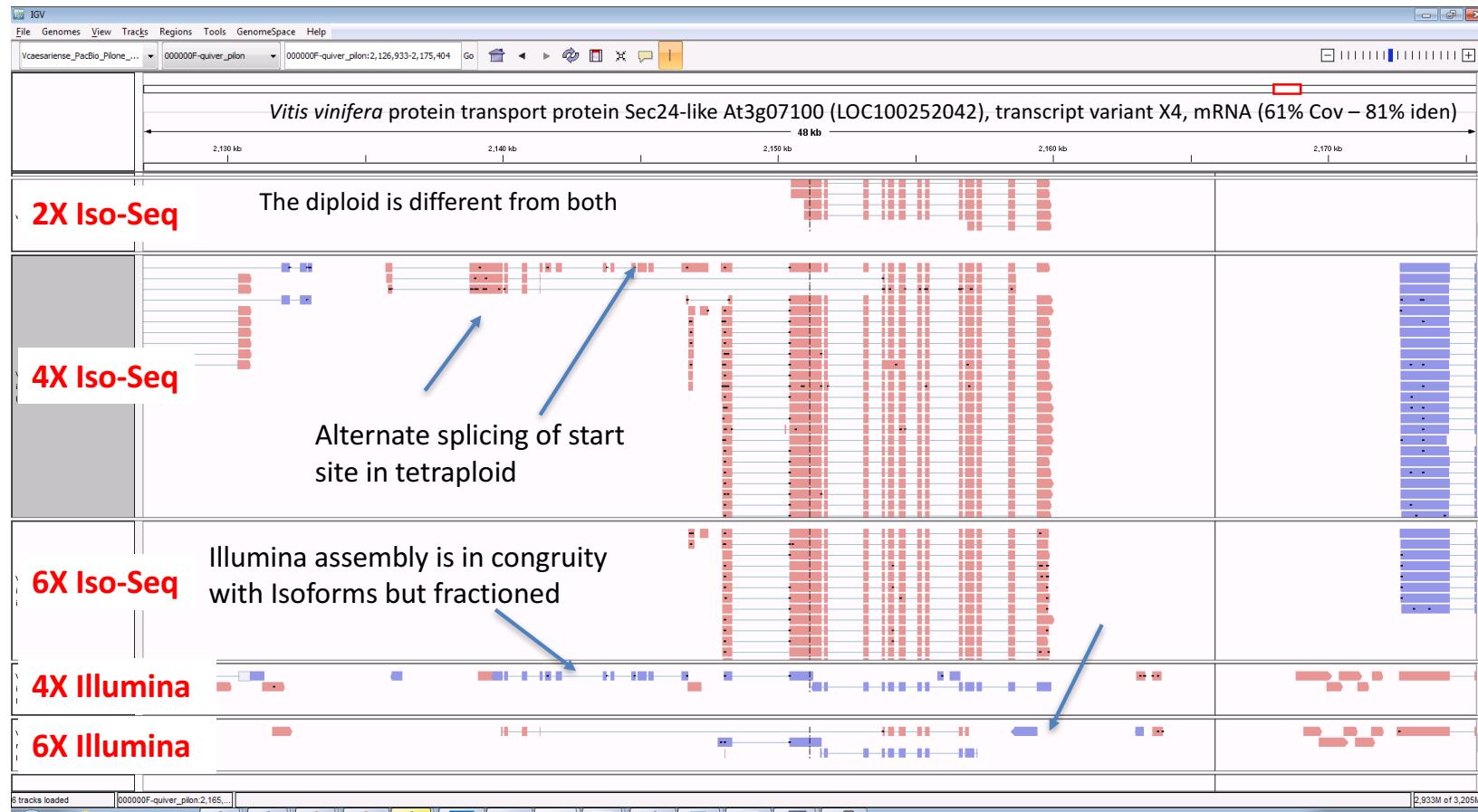
Polymorphisms in hexaploids not in 2X Introns in Illumina assembly



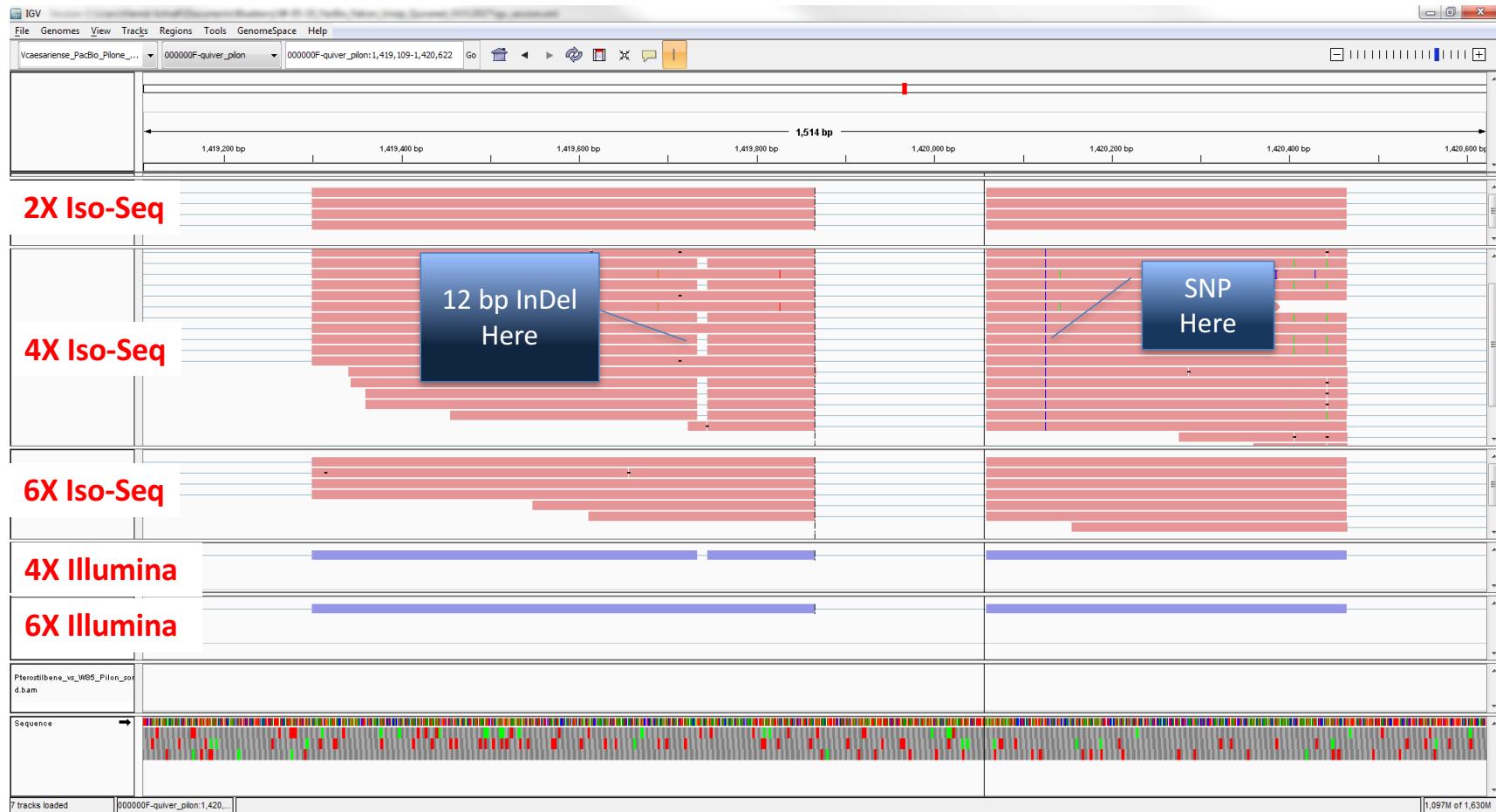




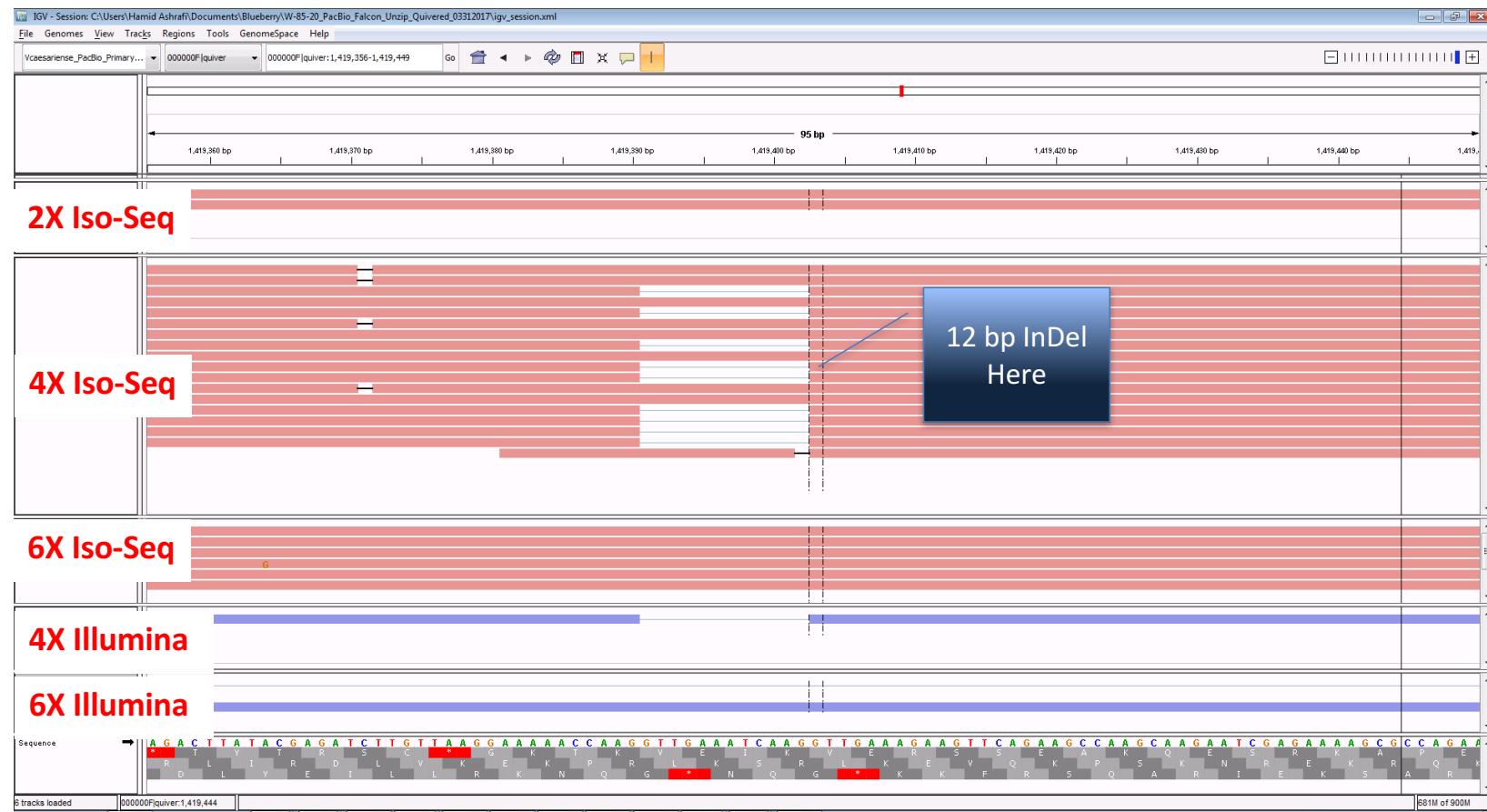
Alternate Splicing of Start Site of a Transport Protein in 4X



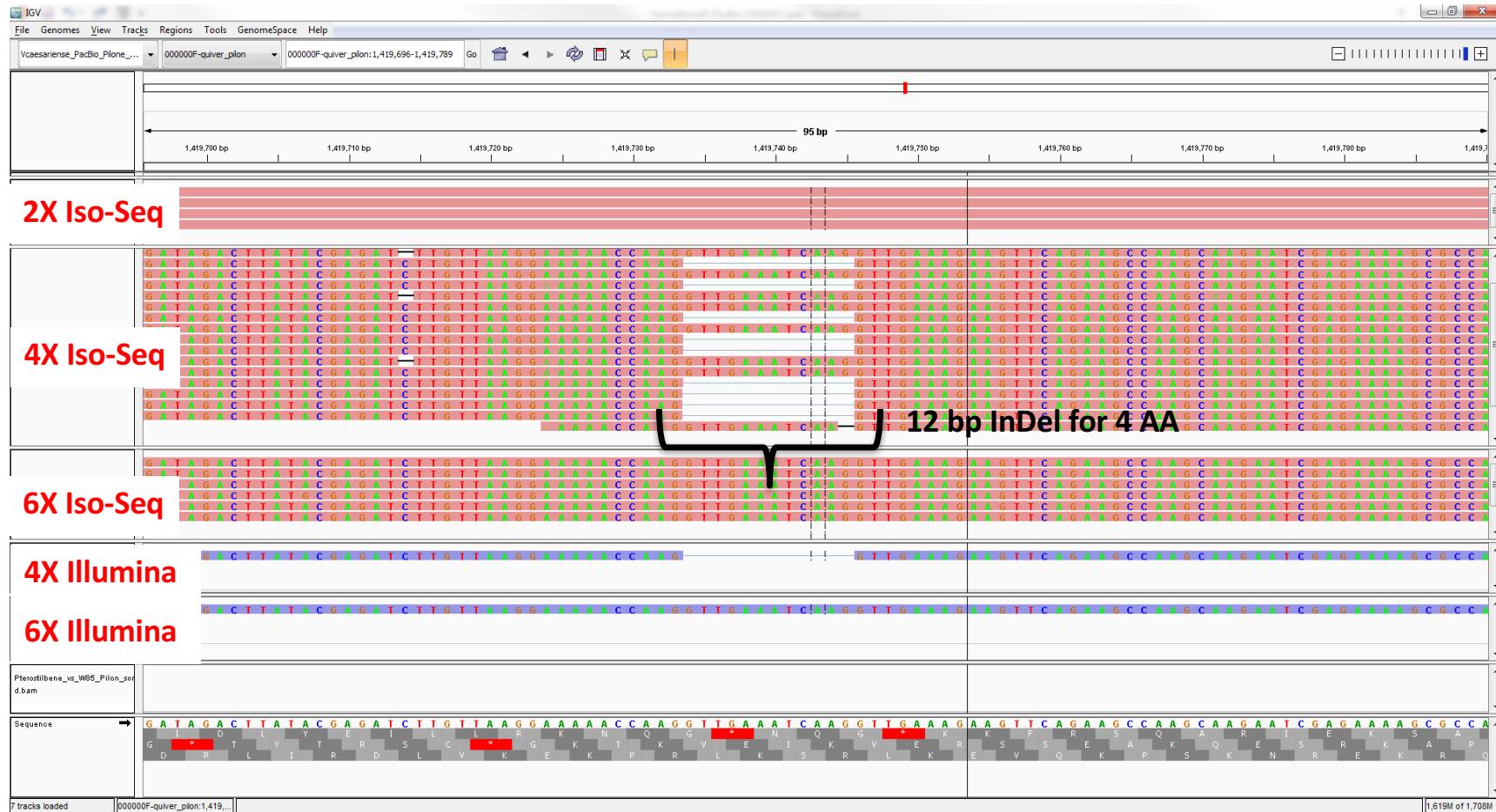
Isoform Analysis Unrevealed Other Forms of Variation (1/3)



Isoform Analysis Unrevealed Other Forms of Variation (2/3)

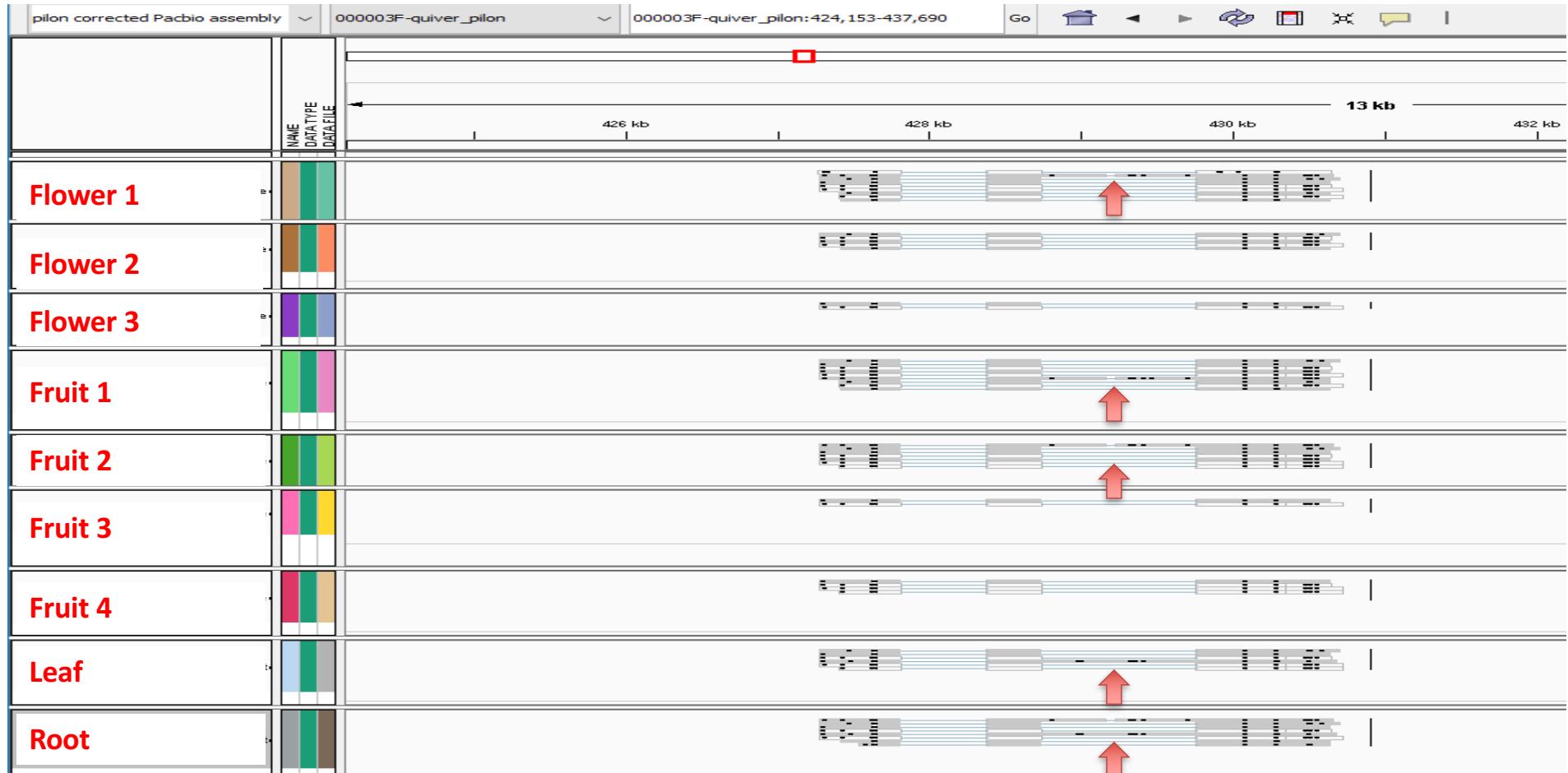


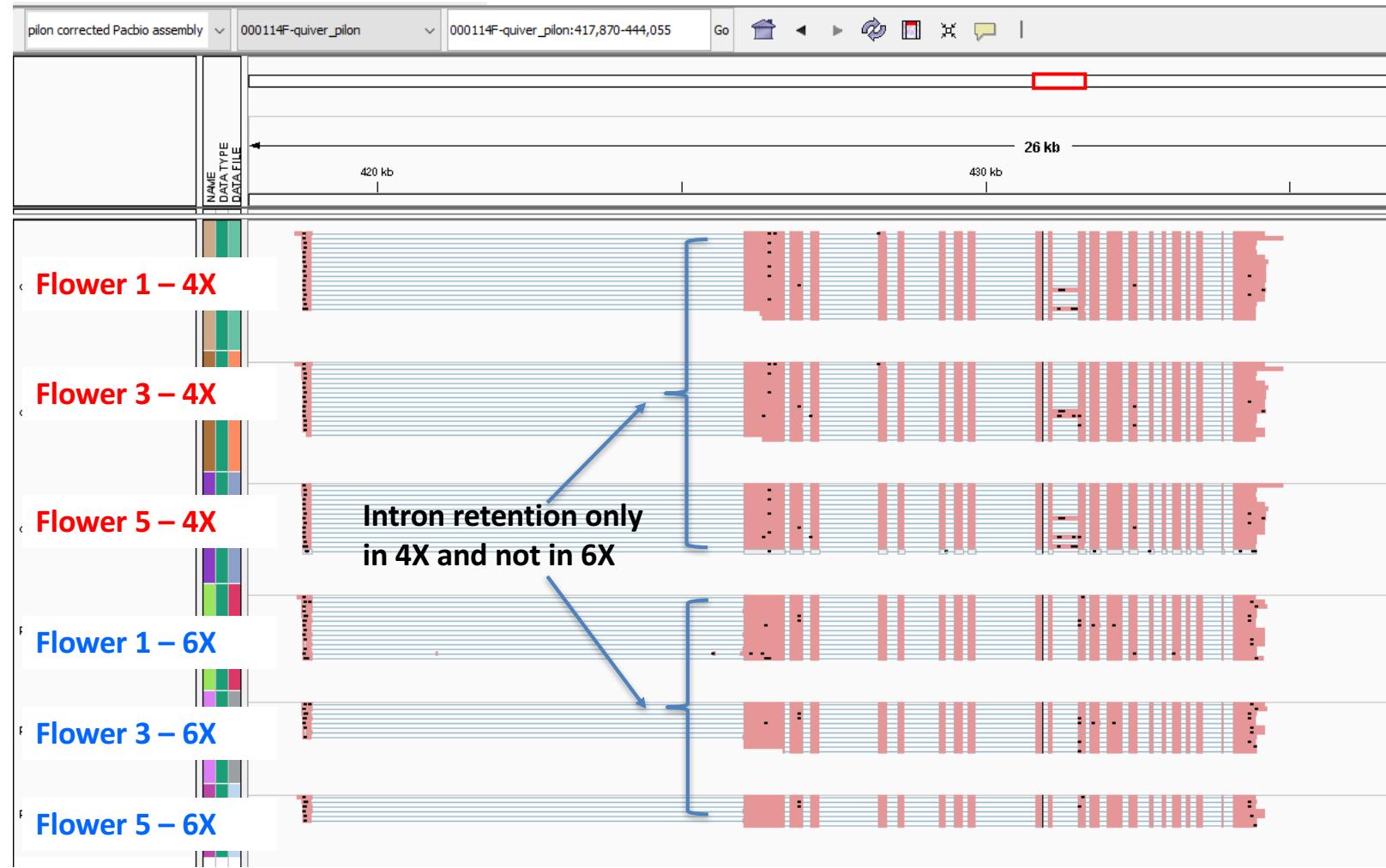
Isoform Analysis Unrevealed Other Forms of Variation (3/3)

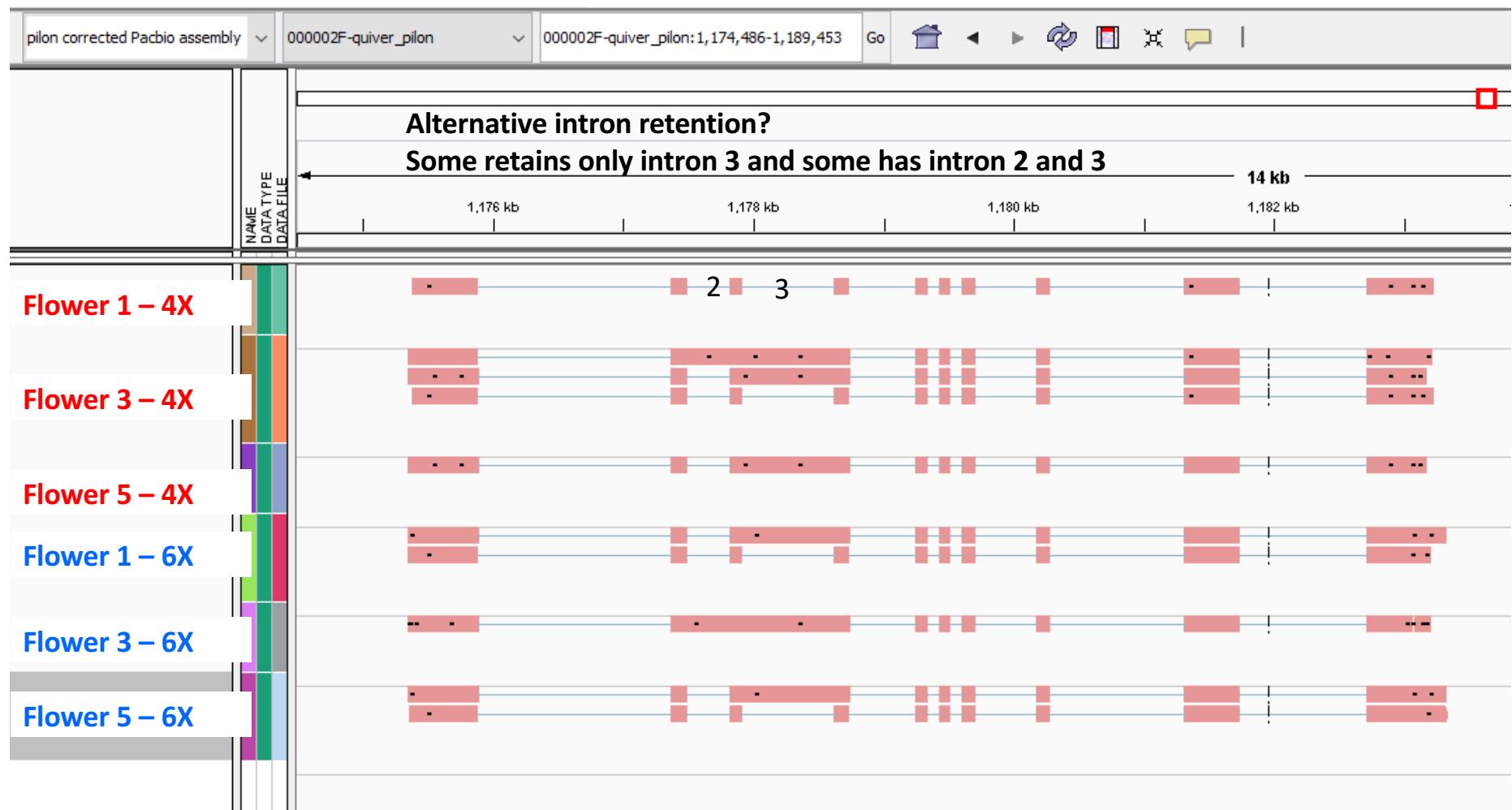


Tissue Specific Iso-Forms

Tissue Specific Intron Retention







Conclusions

- Using 54 SMRT RSII cells we were able to achieve a decent assembly of blueberry genome.
- Barcoding the RNA-Seq libraries helped us to save for the cost of the project.
- What makes different blueberries is more than simple SNP in the genes.
- We believe alternative splicing, Iso-forms and structural variants are responsible for a large proportion of variations and evolution of blueberry.
- We need to use the term “autopolyplody” with caution for blueberry and maybe other plants species.
- Although the chromosomes can pair during meiosis, this does not mean that sub-genomes are identical at the gene and iso-forms levels.
- Defining the iso-forms to have a biological meaning and connecting them to the actual function of the genes remains a challenging task.

**Thank You
Questions?**