

Assembly and Analysis of DNA Sequencing Data

Xiaoqiu Huang

Department of Computer Science

Iowa State University

Raw DNA Sequencing Data are Huge and Short

- A huge number (up to a billion) of DNA sequences can be produced by a sequencing instrument in a day.
- DNA sequences are short (e.g. 100 to 200 base pairs in length).
- Their error rates are often below 1%.

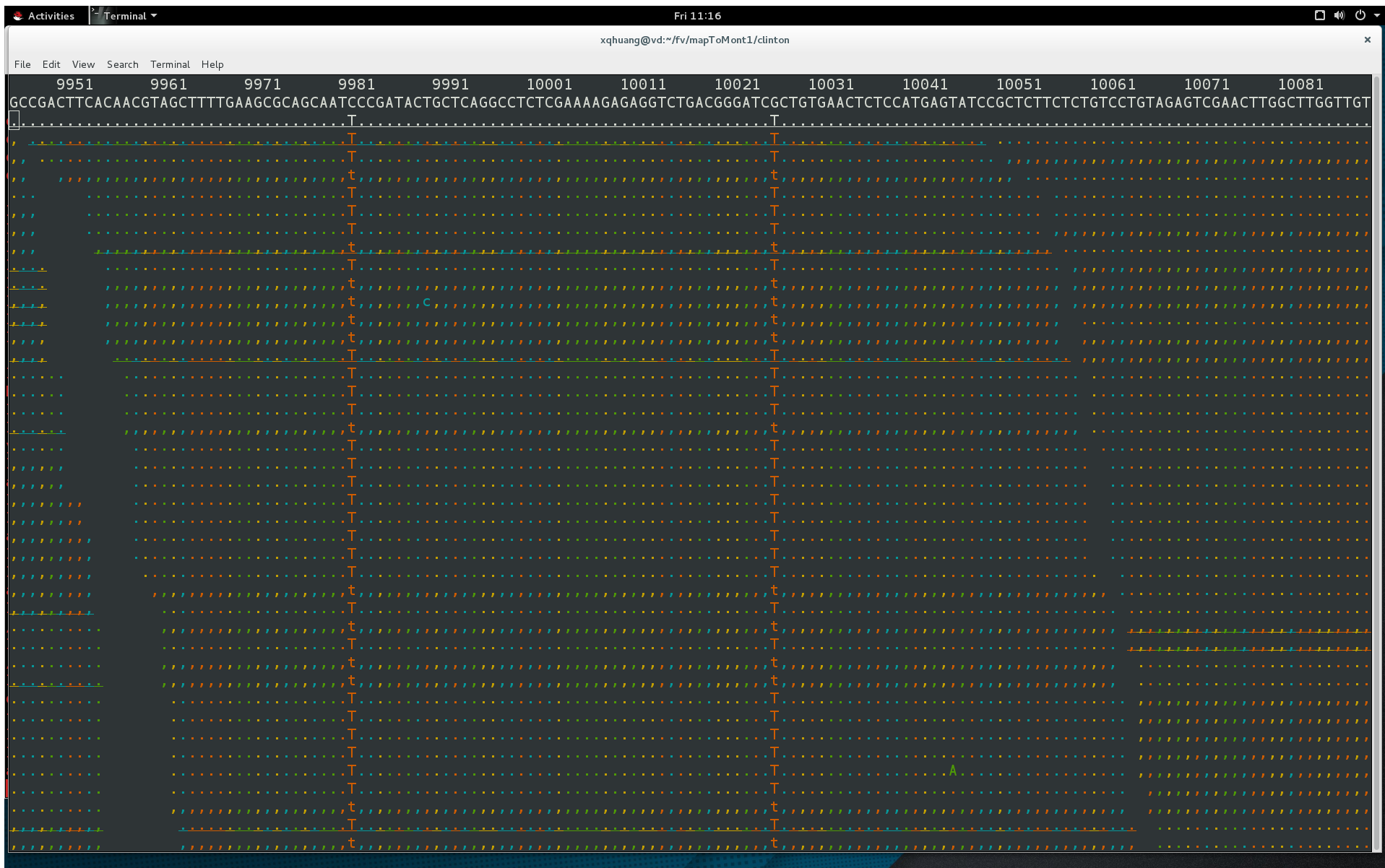
[illegible]

Assembly and Mapping

- Short redundant DNA sequences are put together to construct longer non-redundant genomic or RNA sequences (assembly).
- DNA sequences are aligned to a reference sequence (mapping) to find single nucleotide polymorphisms (SNPs).

A*GAGAGCAAG*TGCCGTGGGCCTGACATGGTCGAGGTGGAAAC*GAGCGTGGGTGA*TGGCGGGCCG*CCACCATATT*CAGTCTTGTCTTAATCC

[illegible]



Methods for Addressing the Challenges

- Dynamic programming
- Hashing
- Graph algorithms
- Parallel processing

Impact of Finding Genetic Variation

- Molecular evolution
- Association of genetic variation with traits
- Personalized lifestyle
- Precision agriculture

Genetic Variation and Reproduction

- Genetic variation is the ‘raw material’ for natural selection.
- Sexual reproduction generates genetic variation by recombination.
- Asexual reproduction lacks mechanisms to generate genetic variation.

Asexual reproduction in Fungi

- Many important plant pathogens have a long period of asexual reproduction.
- Two asexual plant pathogens:

the fungus *Verticillium dahliae*

the fungus *Fusarium virguliforme*

Fungus *Fusarium virguliforme*

- Sudden death syndrome (SDS) of soybean first appeared in Arkansas in 1971 and then spread in 30 years to all major soybean-producing regions in the U.S.
- SDS in North America is mostly caused by the fungus *Fusarium virguliforme* (Fv).
- SDS is similar to Bean Root Rot (BRR).

Low Genetic Diversity but Variable Aggressiveness

- Low genetic diversity is detected in *Fv* populations based on small sequence data.
- But *Fv* isolates show variable aggressiveness on soybean.
- Paradox: *Fv* is subject to Muller's Ratchet but is highly adaptive.

Genome Sequence Data

- 454 sequences of up to 400 bp:
3 Gb of data for *Fv* isolate Mont-1,
single reads and paired reads with 3-kb
and 20-kb inserts.
- Illumina sequences of 102 bp:
390 Gb of data for 10 SDS/BRR isolates,
39 Gb of paired-end reads with a 300-bp
insert per isolate.

Experimental Results

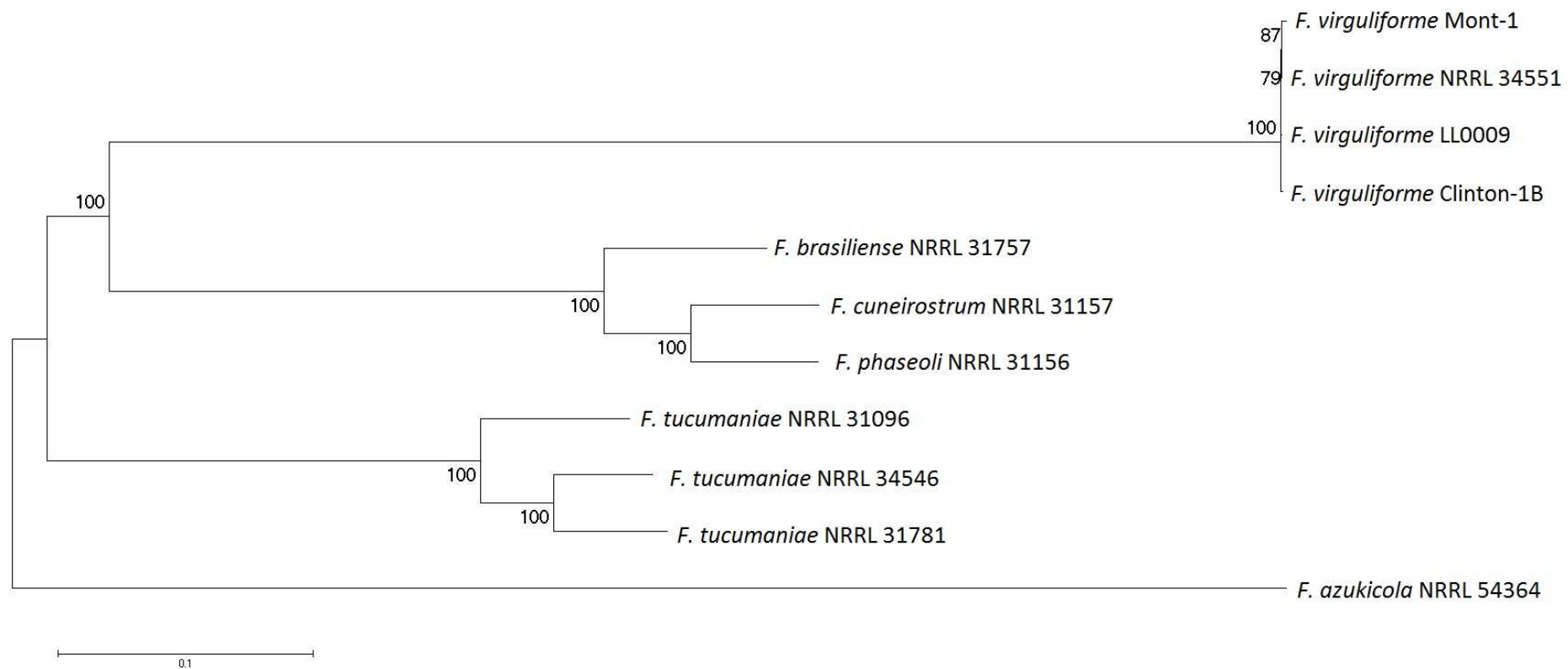
- PCAP was used to produce a chicken genome assembly on ABI 3730 reads (Hillier et al., Nature 432, 695-716).
- PCAP.454 was used to produce a genome assembly of *Fv* Mont-1 (Srivastava et al. PLoS One 9(1), e81832).
- PCAP.Illumina produced an assembly on each of the 10 SDS/BRR isolates.

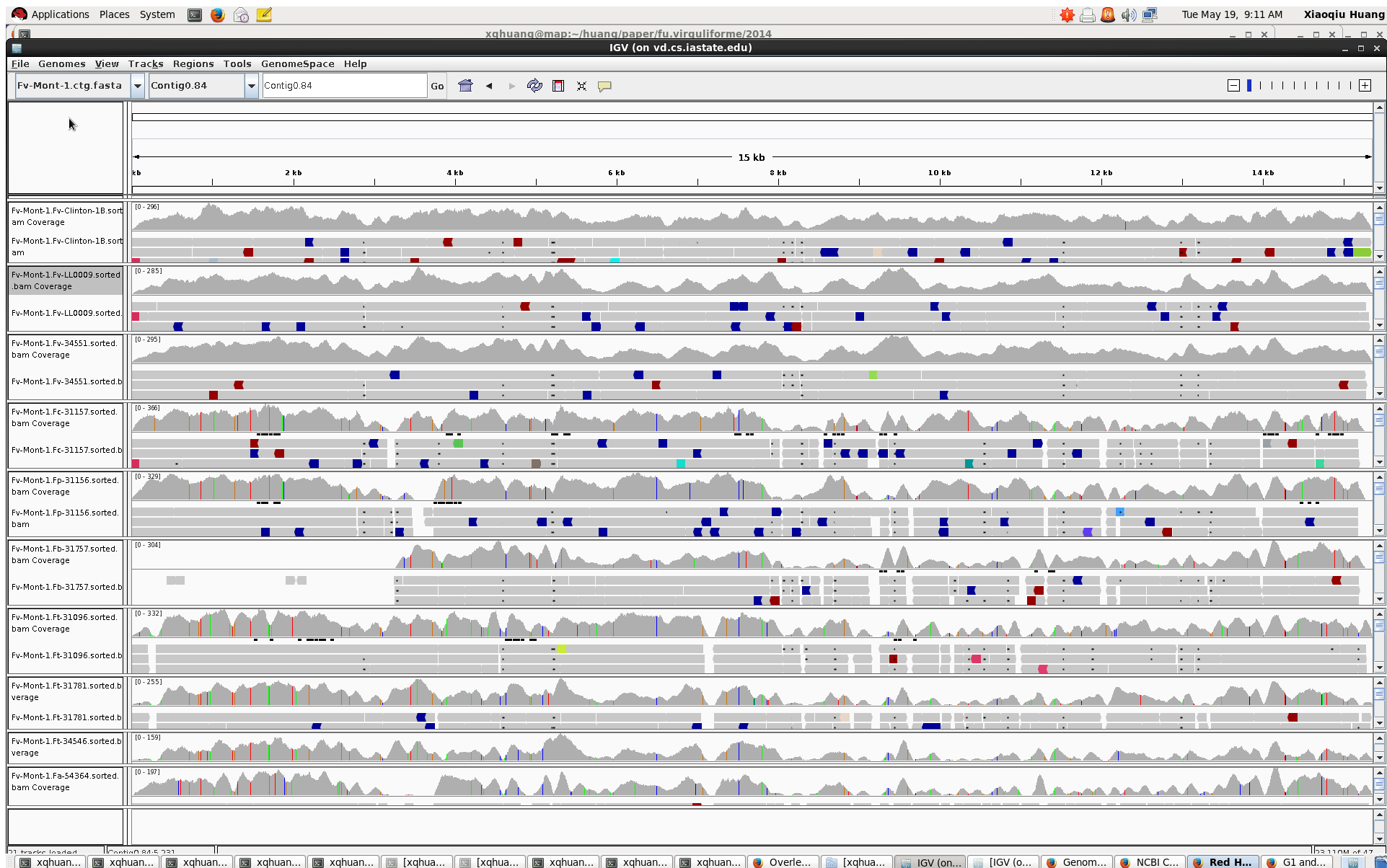
Data Comparison

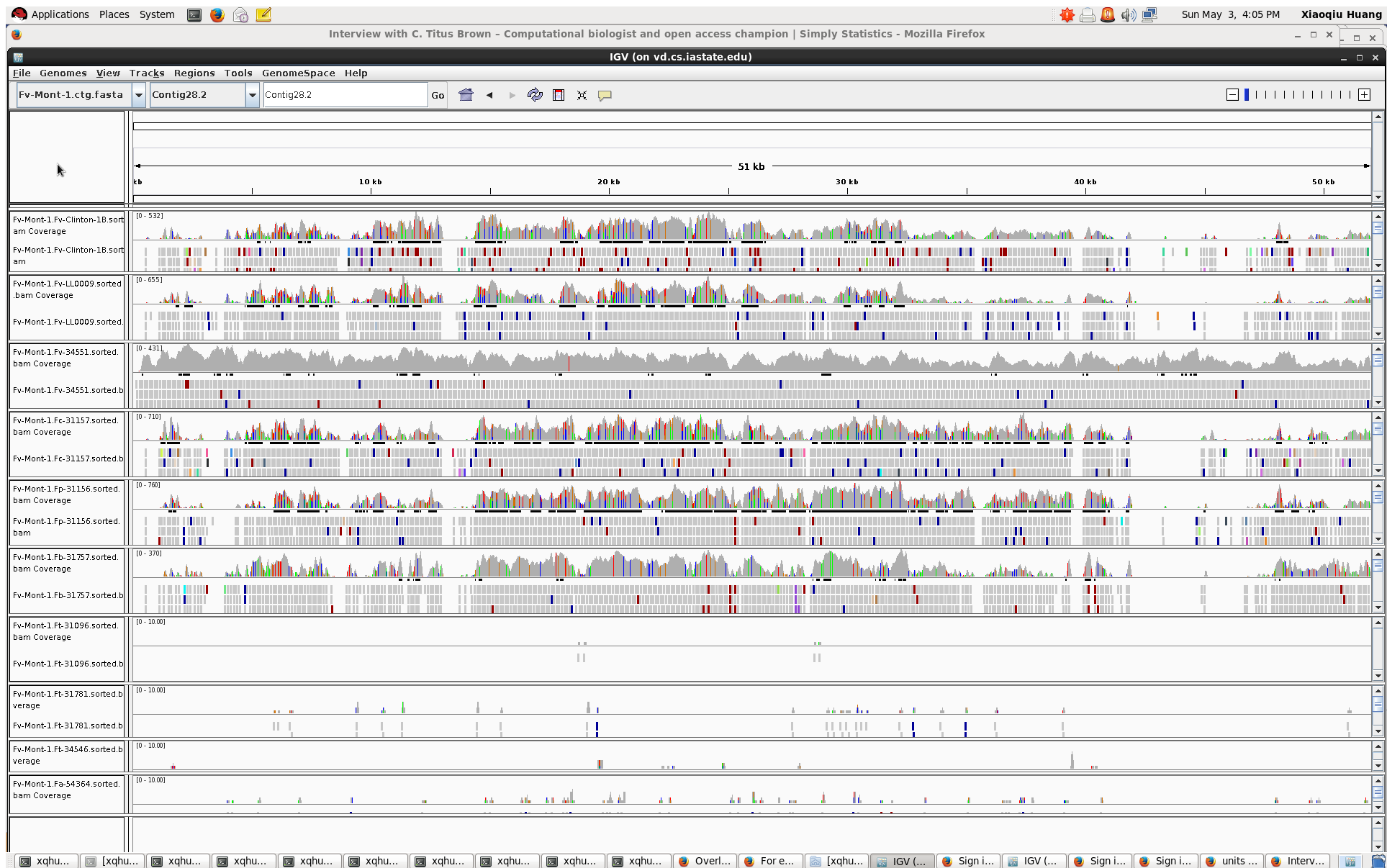
Type	No. of reads (in millions)	Size (Gb)
3730 Chicken	9.8	20
454 Fungus <i>Fv</i>	5	3
Illumina F. <i>Ft</i>	151	44

Computation Time Comparison

Type	#Processors	Time
3730 Chicken	100	1 week
454 Fungus Fv	8	2 days
Illumina F. Ft	1	2 days







Activities

Document Viewer

Fri 14:47

1 of 1

<

>

Q

SNP

SNP.pdf

214.32%

Q

≡

×

Thumbnails

1

Table 1: Length of coverage and distribution of SNPs when reads were mapped onto reference *Fv* Mont-1

Isolate	Length of coverage (Mb)	Number of SNPs	Mean SNP rate/ standard deviation ^a	Max SNP rate ^b
<i>Fv</i> 34551	49.9	4,955	0.00003/0.00007	0.00177 (23.7)
<i>Fv</i> Clinton-1B	49.6	8,269	0.00006/0.00044	0.00960 (21.5)
<i>Fv</i> LL0009	49.2	8,541	0.00007/0.00052	0.01129 (21.7)
<i>Fp</i> 31156	40.0	178,511	0.00447/0.00125	0.01097 (5.2)
<i>Fc</i> 31157	39.5	176,065	0.00446/0.00123	0.01126 (5.5)
<i>Fb</i> 31757	39.3	172,100	0.00435/0.00117	0.00903 (4.0)
<i>Ft</i> 31096	39.3	181,420	0.00462/0.00128	0.00943 (3.8)
<i>Ft</i> 31781	39.2	172,823	0.00441/0.00114	0.00829 (3.4)
<i>Ft</i> 34546	38.9	157,076	0.00412/0.00102	0.00726 (3.1)
<i>Fa</i> 54363	37.9	188,209	0.00506/0.00119	0.00957 (3.8)

^a The mapped reference was partitioned into at least 700 disjoint windows each with 35-kb sufficiently covered base positions. The mean and standard deviation were calculated for the SNP rates of these windows.

^b The number in the parentheses is the maximum SNP rate measured in units of standard deviation above the mean.

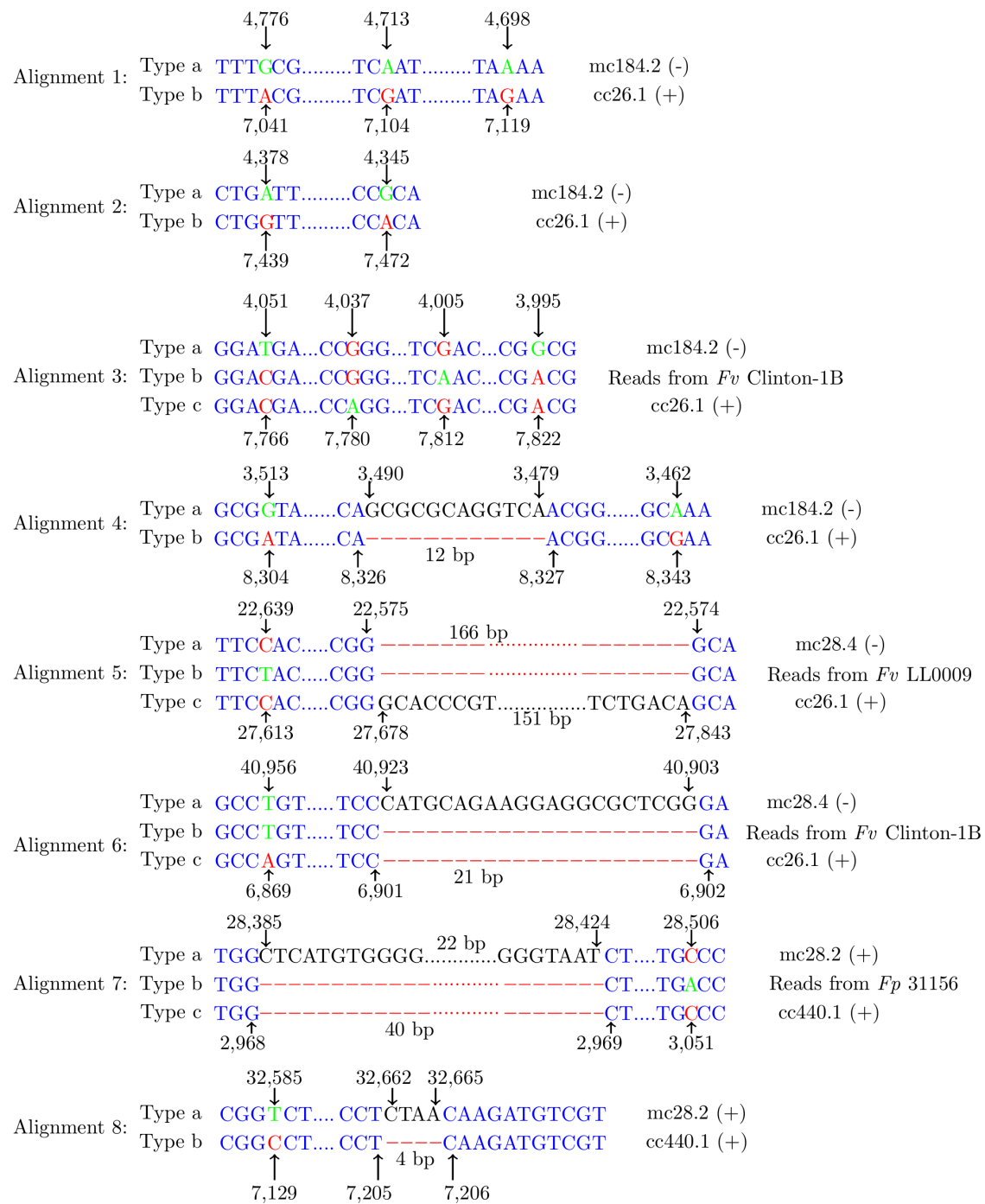


Table 1: The number of reads from the isolate that link all alleles in the sequence

Sequence ^a	Number of reads from the isolate that cover the sequence ^b					
	<i>Fv</i> Clinton-1B	<i>Fv</i> LL0009	<i>Fv</i> 34551	<i>Fc</i> 31157	<i>Fp</i> 31156	<i>Fb</i> 31757
A1.Ta	16	16	16	0	0	0
A1.Tb	32	18	0	48	46	8
A2.Ta	88	149	147	0	0	9
A2.Tb	85	114	0	115	0	0
A3.Ta	52	41	61	0	0	0
A3.Tb	78	52	0	152	0	0
A3.Tc	33	0	0	34	46	0
A4.Ta	162	134	77	121	0	97
A4.Tb	54	0	0	127	242	0
A5.Ta	39	27	57	0	0	18
A5.Tb	0	8	0	39	0	0
A5.Tc	85	0	0	69	65	0
A6.Ta	0	0	46	0	0	0
A6.Tb	72	0	0	0	209	0
A6.Tc	116	121	35	554	244	74
A7.Ta	0	0	98	0	0	0
A7.Tb	0	31	0	0	54	0
A7.Tc	35	0	0	42	50	39
A8.Ta	0	35	34	40	42	0
A8.Tb	70	0	0	51	59	42

^a Each sequence is denoted by its alignment number and type letter (Figure X): e.g., Types a and b in Alignment 1 are denoted by A1.Ta and A1.Tb, respectively.

^b A read covers a sequence in a set of polymorphic sequences if the read and the sequence have the same allele at every occurrence of polymorphism.

Conclusion

- Low SNP rates show that the assembly consensus sequence is accurate.
- They confirm that *Fv* is asexual in reproduction mode.
- *Fv* genome has mutational hot regions.
- *Fv* has novel mechanisms for generating mutations.

Acknowledgements

- Madan Bhattacharyya at ISU
- Leonor Leandro at ISU
- Kerry O' Donnell at USDA ARS
- Subodh Srivastava at U. of Arkansas
- Shiaw-Pyng Yang at Monsanto
- Colleagues at WashU