



Manual Sequence Improvement of Bacterial Genomes

Aye Wollam

Genome Sequencing Center
Washington University School of Medicine
St. Louis, MO

awollam@watson.wustl.edu

Purpose

- To produce sequence assemblies that are qualitatively better than draft assemblies.
- To do so in less cost, time and labor than those needed for finishing genomes.

awollam@watson.wustl.edu

How we are using it

- For investigative sampling of a population, such as Human Gut Microbiome Project.
- 100 bacterial genomes of human gut are being sequenced, manually improved, and annotated.
- Selected genomes will be finished to the full.
- 1000 genomes in future pipeline

awollam@watson.wustl.edu

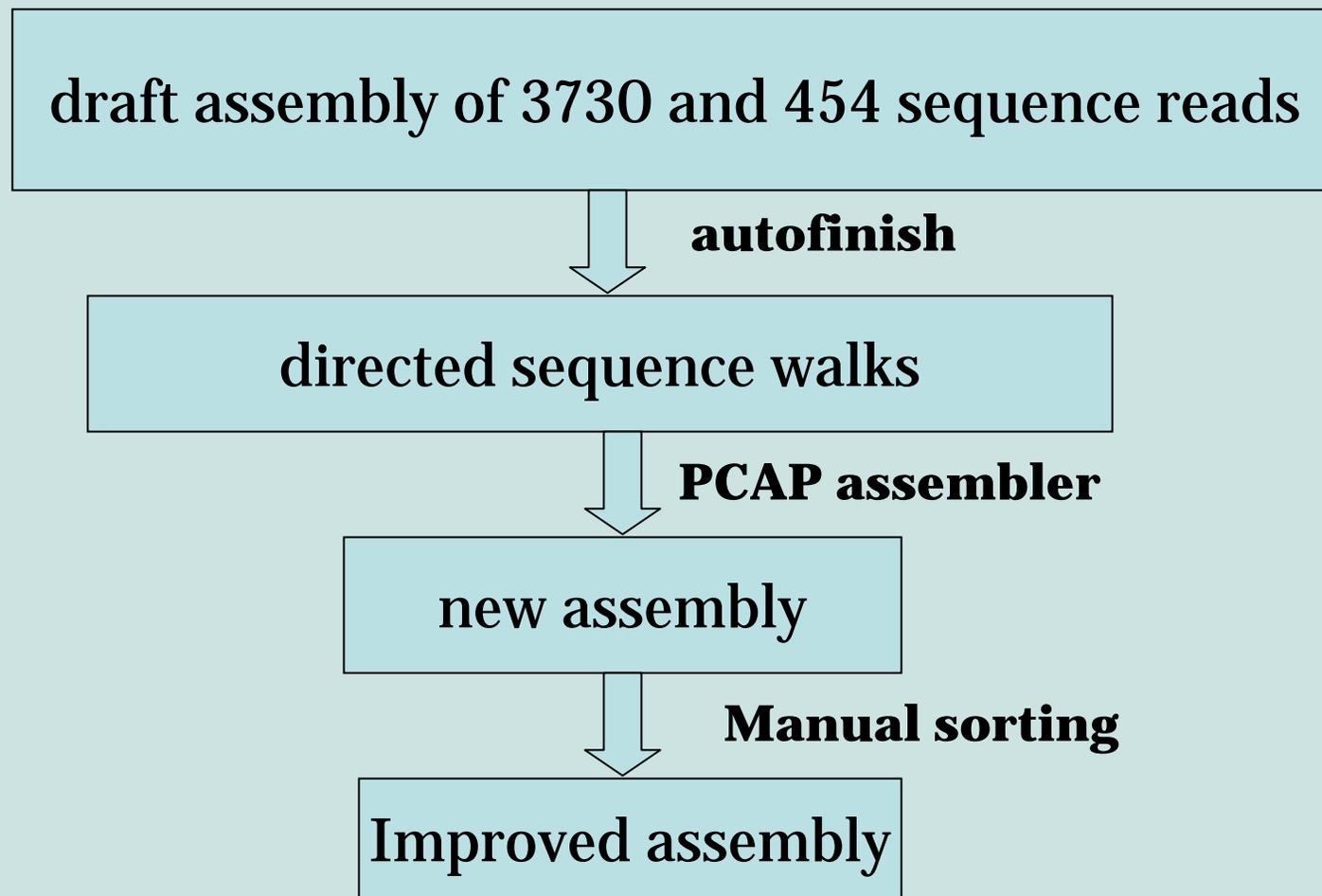


What are the steps in manual
sequence improvement of
bacterial genomes?

awollam@watson.wustl.edu



Steps involved in manual sequence improvement



Manual sorting process

- Perform joins missed by assembler
- Correct consensus errors
- Sort misplaced reads
- Identify and correct misassembled regions
- Order and orientation of contigs are provided if possible

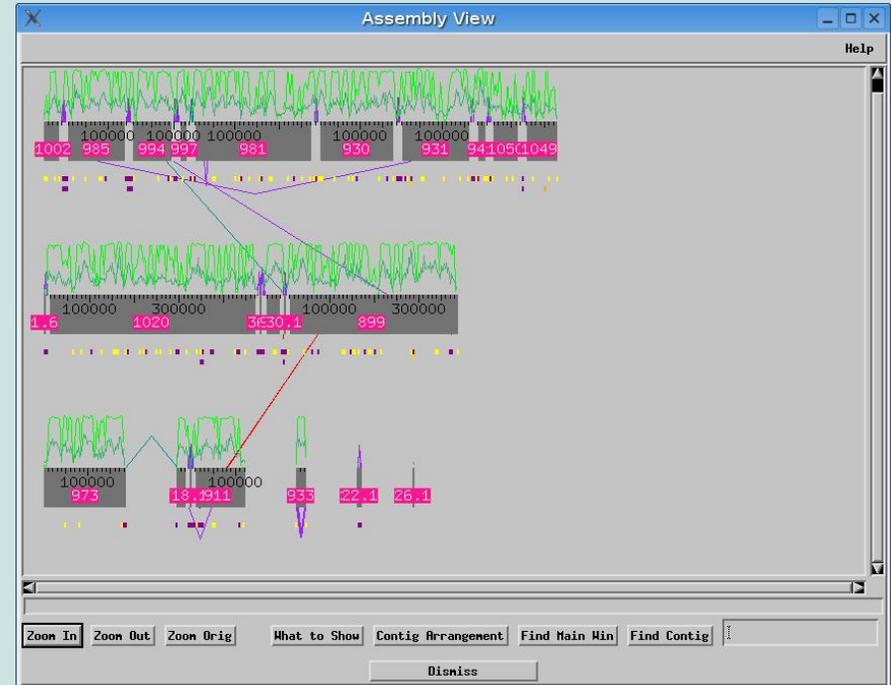
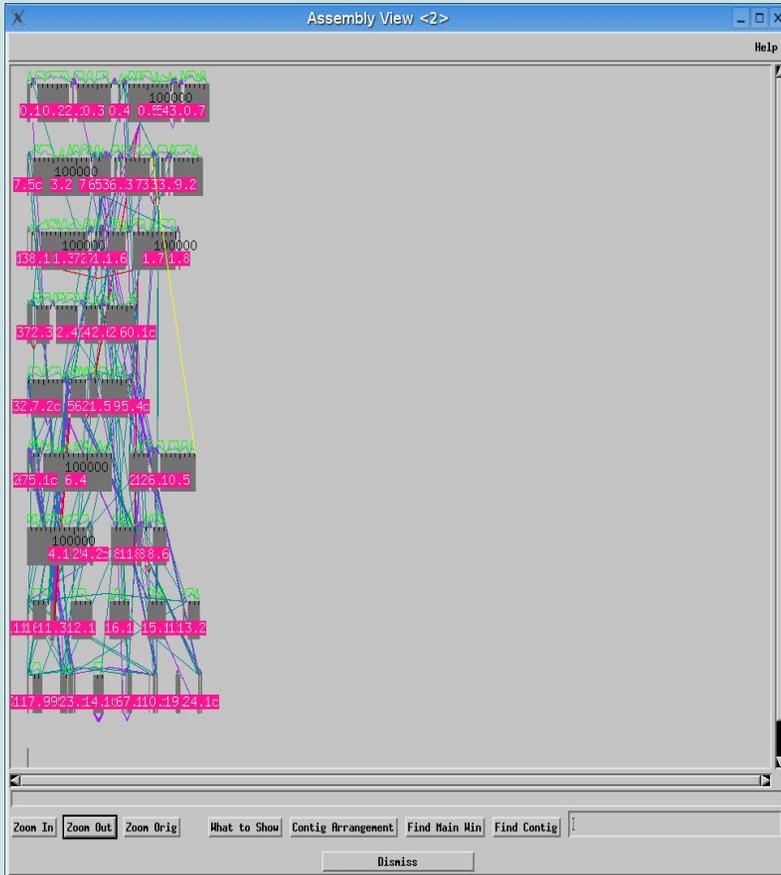
awollam@watson.wustl.edu

Perform joins missed by assembler

- look for missed joins
- sequence searches within a scaffold
- detect overlaps between scaffolds

awollam@watson.wustl.edu

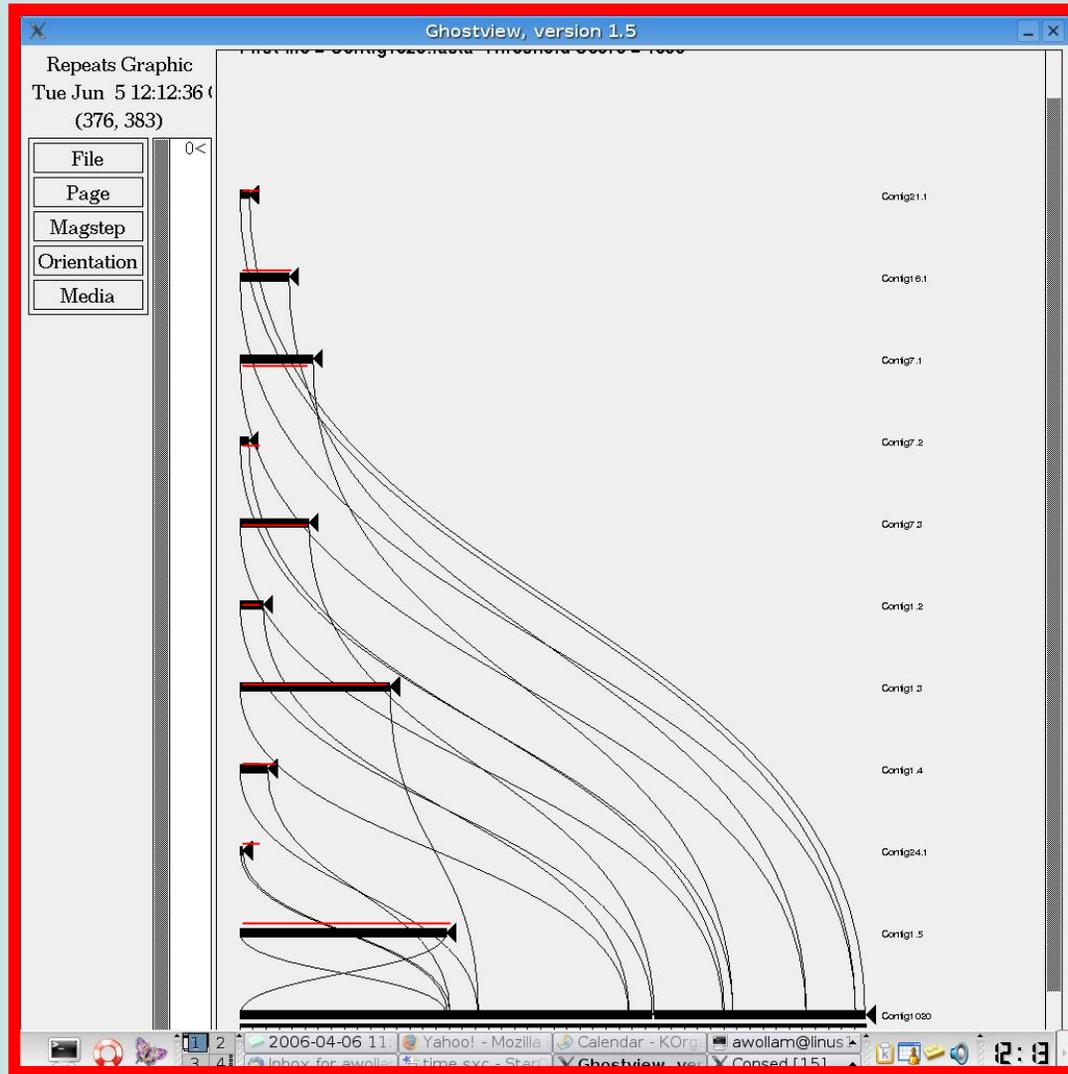
Perform joins missed by assembler



joins not made by assembler;

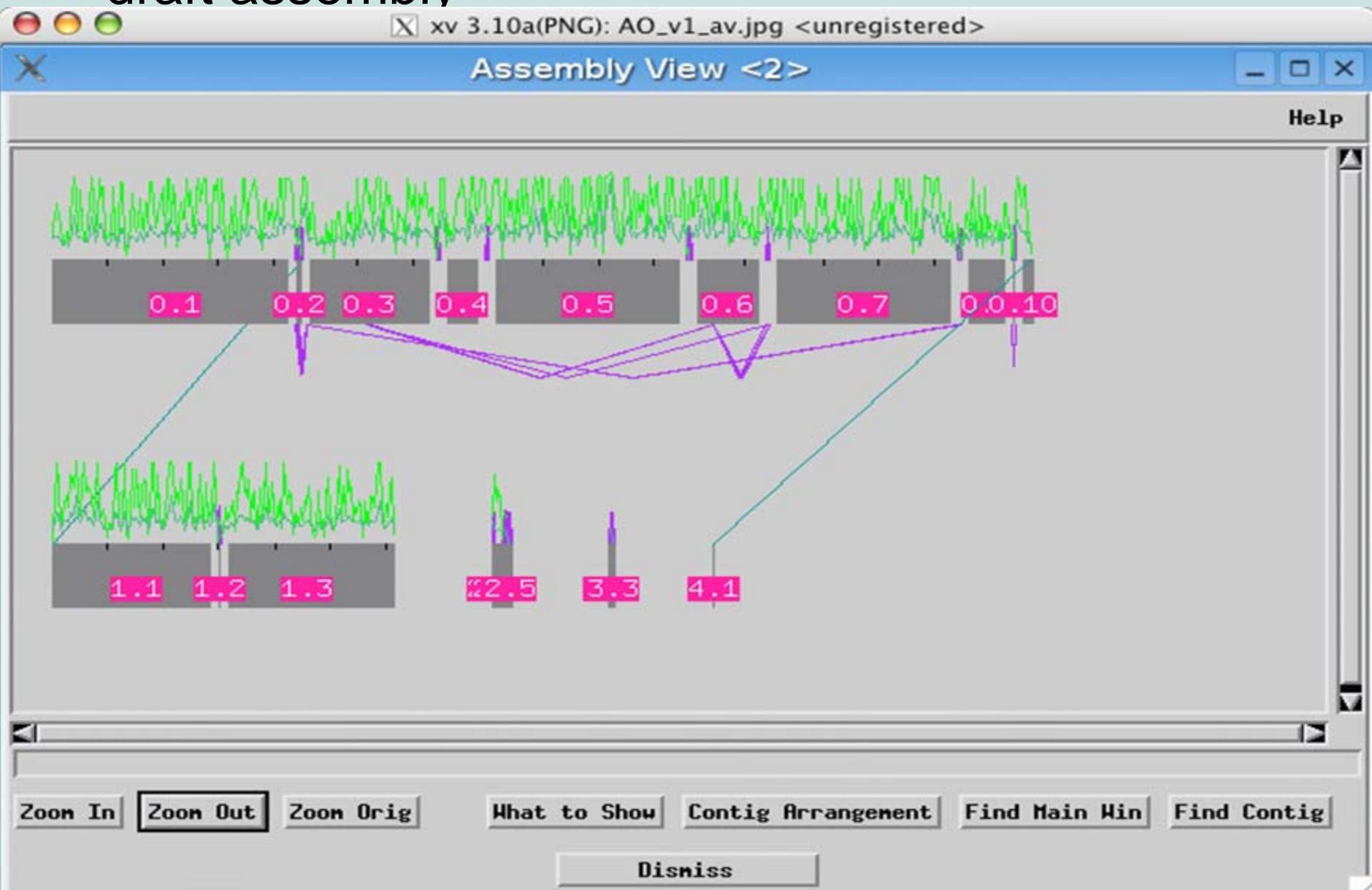
contigs sorted and joins made manually

Perform joins missed by assembler

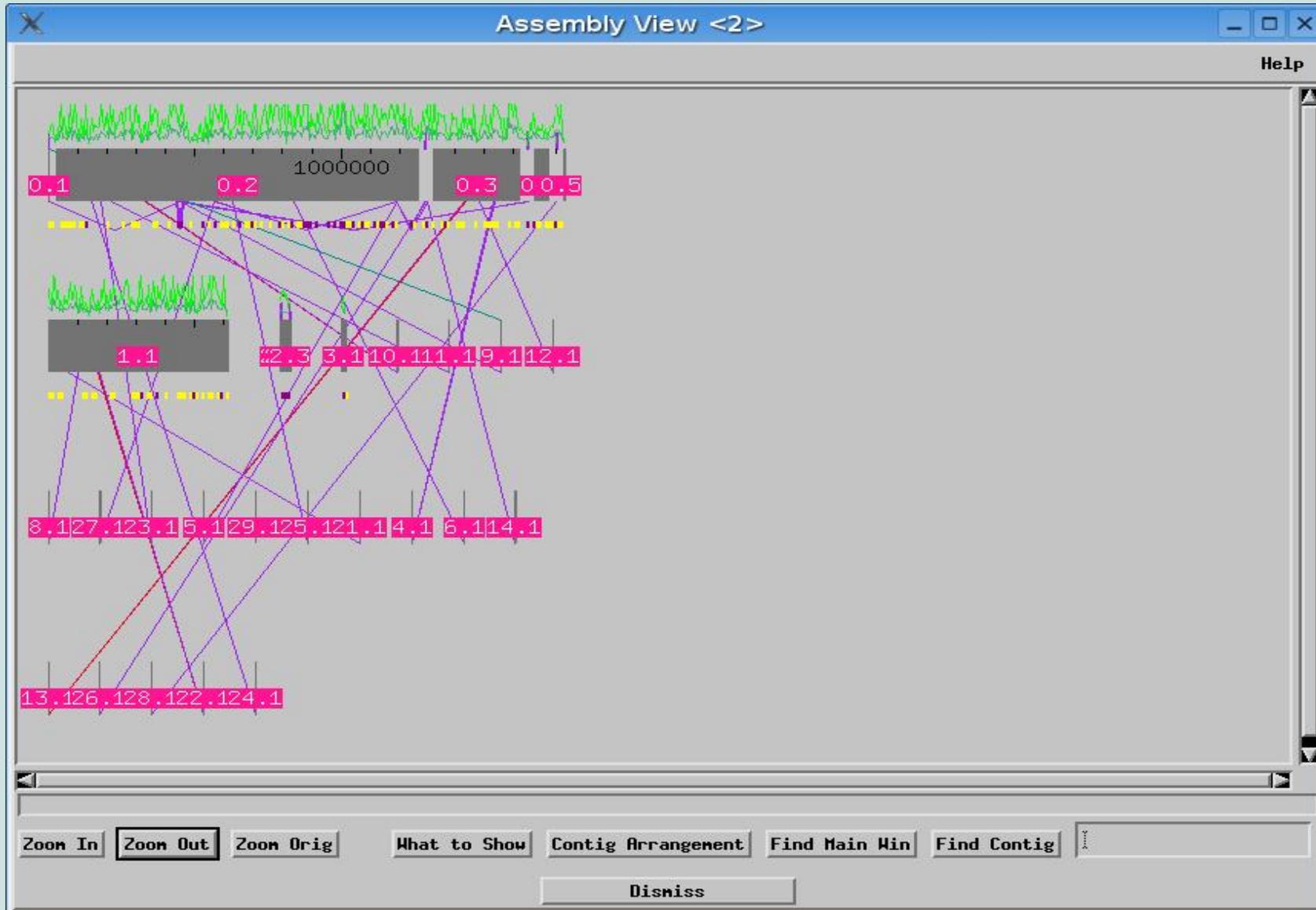


joins made manually in msi process;
red lined contigs(draft contigs) compared with MSI contig

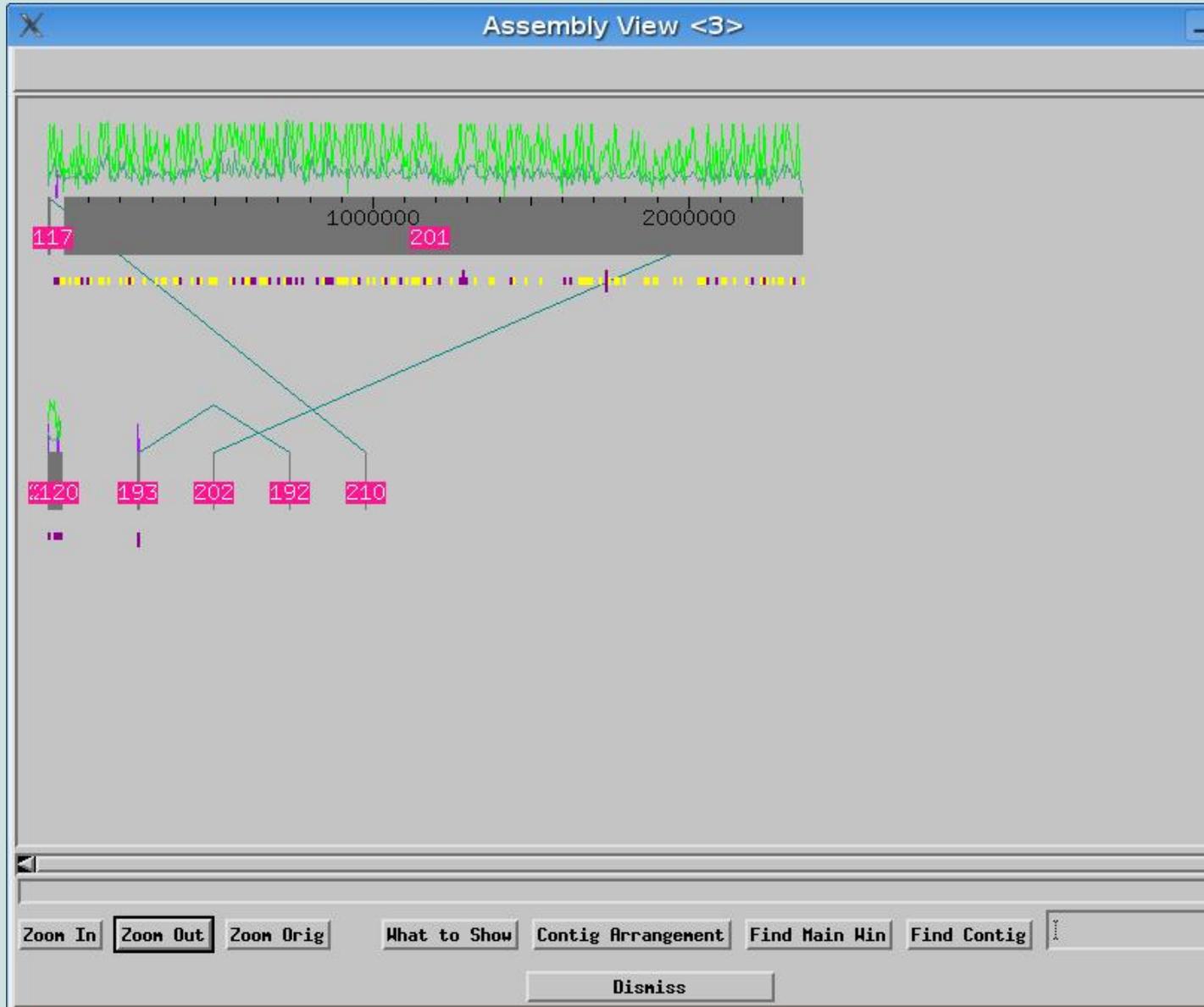
draft assembly



2nd assembly with autofinish data



After manual improvement



watson.wustl.edu

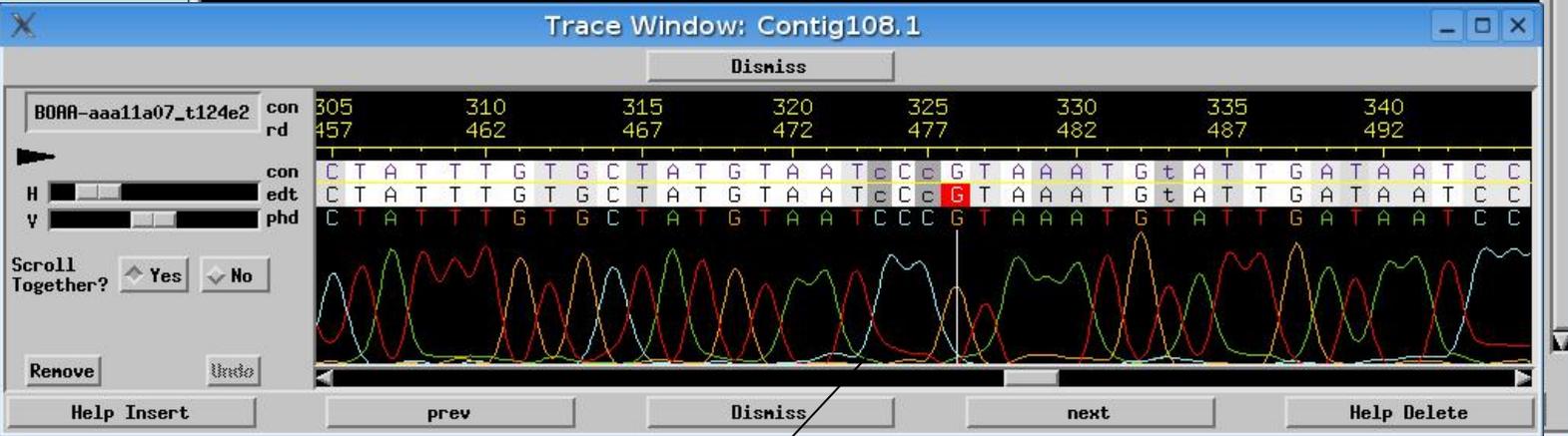
Correct sequence errors

- Navigate by low quality regions and high quality discrepant bases
- correct consensus errors
- trim off irrelevant sequences at contig ends

awollam@watson.wustl.edu



Correct sequence errors



wrong base calls (3 Cs) instead of 2 Cs here

awollam@watson.wustl.edu



Sort misplaced reads

example 1a: sequence walks disconnected to the priming site

The image displays two windows from a genome browser showing aligned reads. The top window, titled "Aligned Reads <17>", shows a contig with a gap. The bottom window, titled "Aligned Reads <16>", shows a read that is disconnected from the priming site. A red arrow points from a box in the top window to a box in the bottom window. A text box explains that the walk should be pulled and connected to the priming site.

Aligned Reads <17>
Bacteroides_uniformis-v3.0.1.fasta.screen.ace.nsi.aw.3 Contig256
Search for String Uncompl Compare Cont Find Main Win Err/10kb: 2541
240 250 260 270 280 290 300 310
CONSENSUS agtaccgggaaaaaatgataagtggaaAATCCagtgtgaagctgc*CCCTAAAACCAAGCGATTCTGCCCCCTT*GTGCTAAA
BUAA-aaa28f03_t187e337.b1 aGTAACGGGAAAaaatgatagttggaaatccagtgtaagctgc*cctAAAACcaagcgattcgcct*cctc*g*gctaAa
BUAA-aab45c10_t186e336.b1 GTTAAACCTCTTTGATAAGCTGgAATAATCCagtgtgaagctgc*CCCTAAAACCAAGCGATTctgocdCttagtgcTAAA

Aligned Reads <16>
Bacteroides_uniformis-v3.0.1.fasta.screen.ace.nsi.aw.3 Contig253
Search for String Conpl Cont Compare Cont Find Main Win Err/10kb: 579.06
1,900 1,910 1,920 1,930 1,940 1,950 1,960 1,970 1,980
CONSENSUS TCCGAAGCTGGAAGATAGACGGACGCACATGGCGGCTATTATTGCTAACGCAGTGAGTGCAGCGT TTAACCTCTTTCTGATAAGCT
superread_1673_4.a1 TCCGAAGCTGGAAGATAAGACGGACGCACATGGCGGCTATTATTGCTAACGCAGTGAGTGCAGCGT TTAACCTCTTTCTGATAAGCT
superread_1673_1.c1 TCCGAAGCTGGAAGATAAGACGGACGCACATGGCGGCTATTATTGCTAACGCAGTGAGTGCAGCGT TTAACCTCTTTCTGATAAGCT

mislplaced walk here should be pulled

This is where the walk should connect and extend into the gap.

awollam@watson.wustl.edu

Sort misplaced reads

example 1b: sequence walks placed to the priming site. Extended 850 base pair into the gap

Aligned Reads

File Navigate Info Color Dim Misc Help

Bacteroides_uniformis-v3.0.1.fasta.screen.ace.msi.aw.4 Contig643 Sone Tags Pos: clear

Search for String Uncompl Compare Cont Find Main Min Err/10kb: 579,13

1,890 1,900 1,910 1,920 1,930 1,940 1,950 1,960 1,970 1,980 1,990

CONSENSUS CCTAATCCGAAGCTGGAAGATAAGACGGACGACATGGCCGCTATTATTGCTAACGC...
superread_1673_3.a1 CCTAATCCGAAGCTGGAAGATAAGACGGACGACATGGCCGCTATTATTGCTAACGCAGTGCAGCCGTTAAACCTCTT
superread_1673_4.a1 CCTAATCCGAAGCTGGAAGATAAGACGGACGACATGGCCGCTATTATTGCTAACGCAGTGCAGCCGTTAAACCTCTTCTGATAAGCT
BUAA-aaa28f03_t188e338.b1 caaccggaataicbgacdaagacgttgaaggggtagggcatgaaacctgaaacatctccacttatcttcagtaacgggaaaaaatgatagttg*aaa**cgg*gttaagg
BUAA-aaa28f03_t187e337.b1 ctcttccatgaa*ccctctacatctcca*ttagcTTcaGTACGGGAAAAaatgatagttggaaaaatccagtgtaagct
BUAA-aab45c10_t186e336.b1 agtactattattgcTAcGCAGTGCAGCCGTTAAACCTCTTCTGATAAGCTGgAAAATCCagtgtaagct
superread_1673_1.c1 CCTAATCCGAAGCTGGAAGATAAGACGGACGACATGGCCGCTATTATTGCTAACGCAGTGCAGCCGTTAAACCTCTTCTGATAAGCT

<<< < < Prev Next > >> >>> cursor dismiss

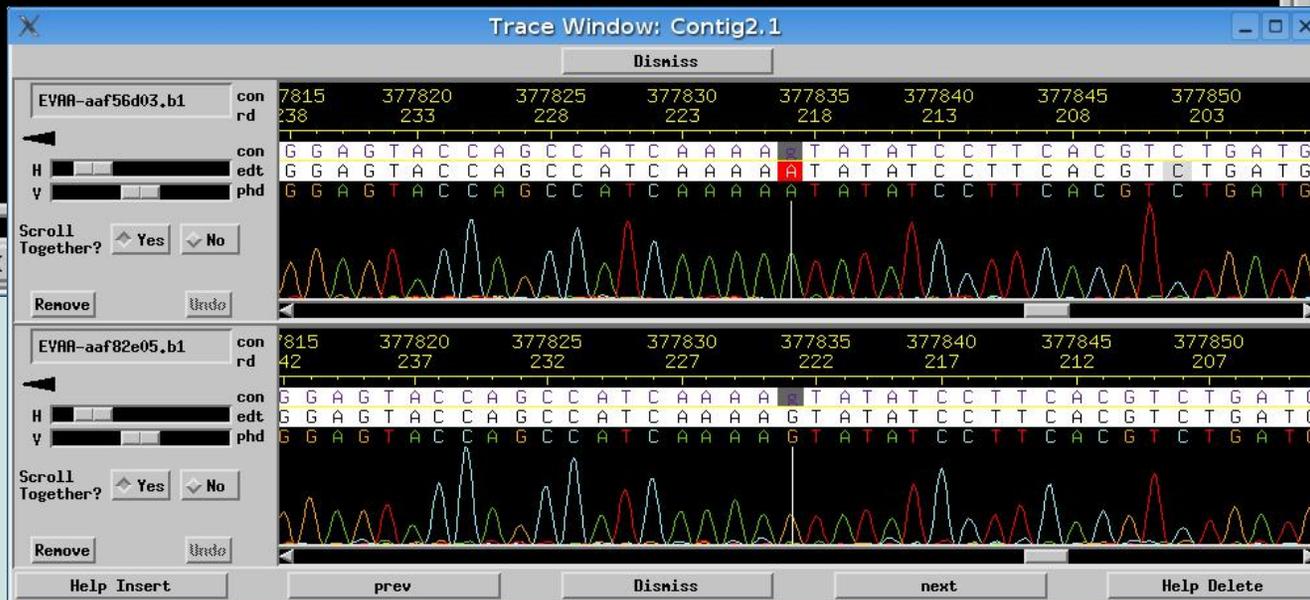
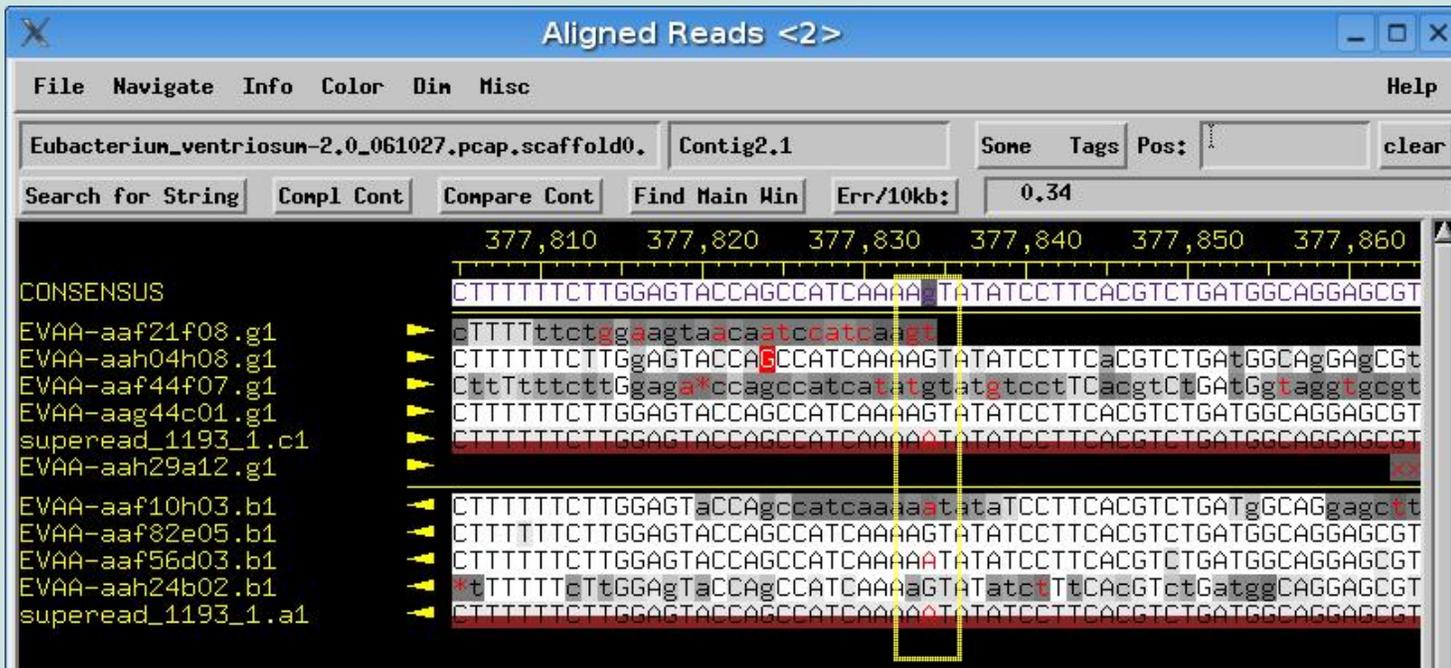
awollam@watson.wustl.edu

Identify and correct misassembled regions

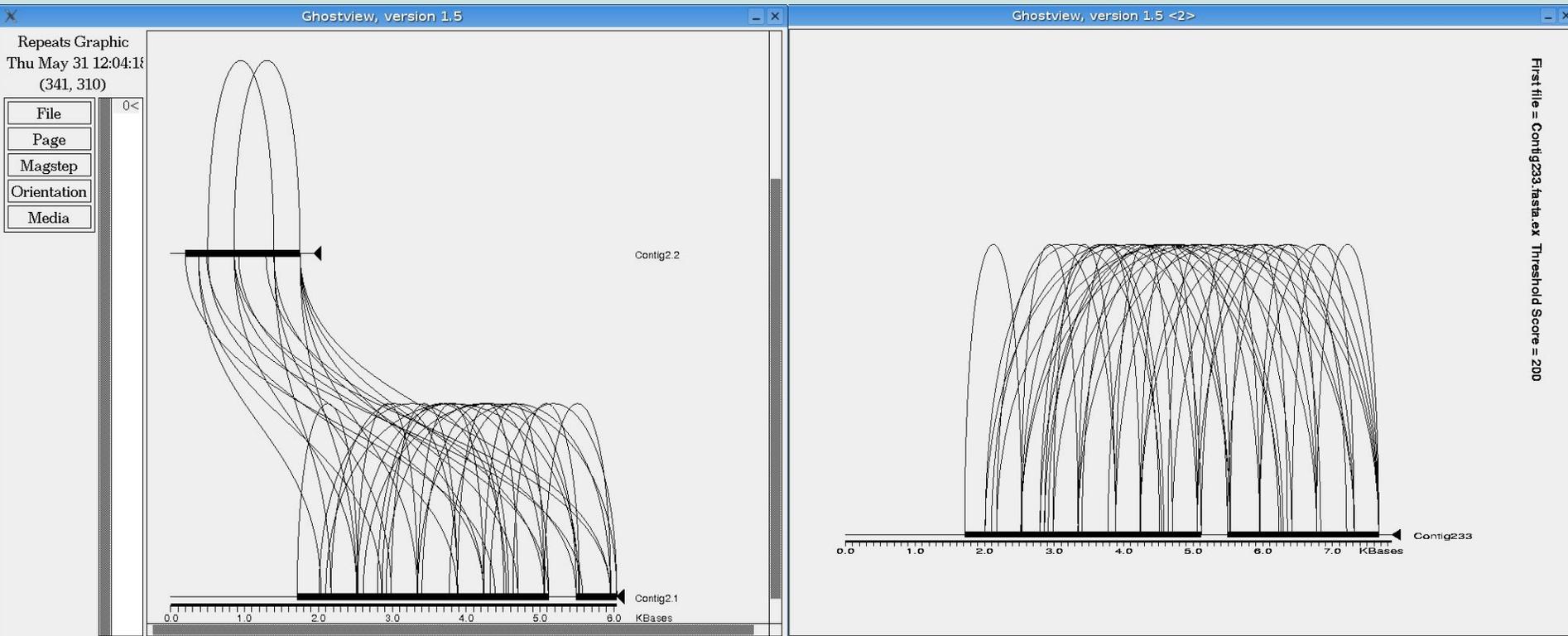
- look for evidence of misassemblies based on read pair inconsistencies and high quality discrepant reads
- rearrange read placement, tear and join contigs as needed.

awollam@watson.wustl.edu

high quality discrepant bases indicating collapsed repeat



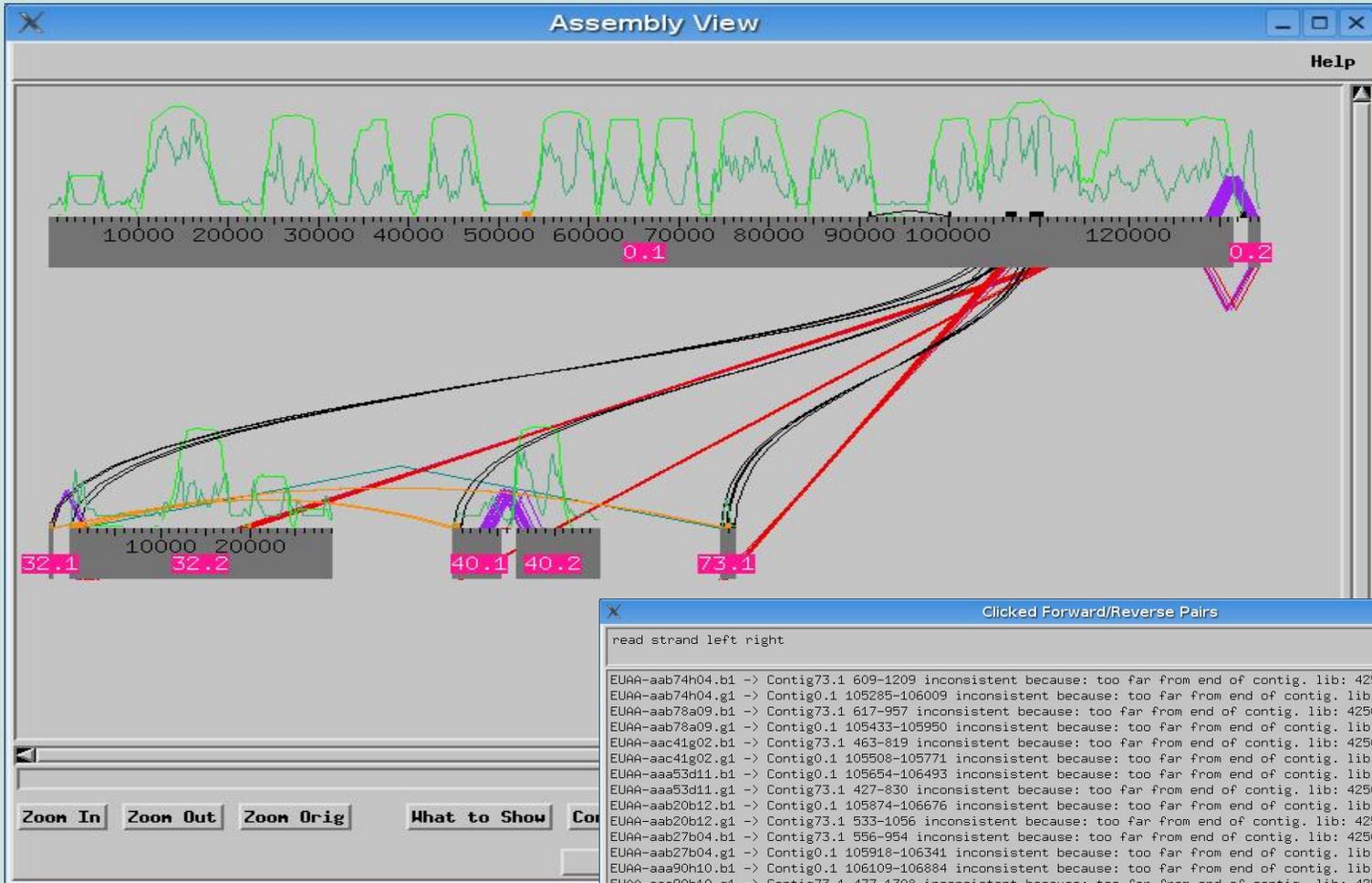
Sort and correct misassembled regions



- two contigs with collapsed repeats are sorted and the result is one contig with repeats separated after manual improvement. No high quality mismatched bases in this contig.

awollam@watson.wustl.edu

Sort and correct misassembled regions



inconsistent read pairs

Sort and correct misassembled regions



- read pairs are now consistent
- repeat copies separated

awollam@watson.wustl.edu

What we don't do in MSI

- additional directed walks \$\$
- PCR for unoriented contigs \$\$
- mini library (short insert libraries) for hard to close gaps \$\$\$

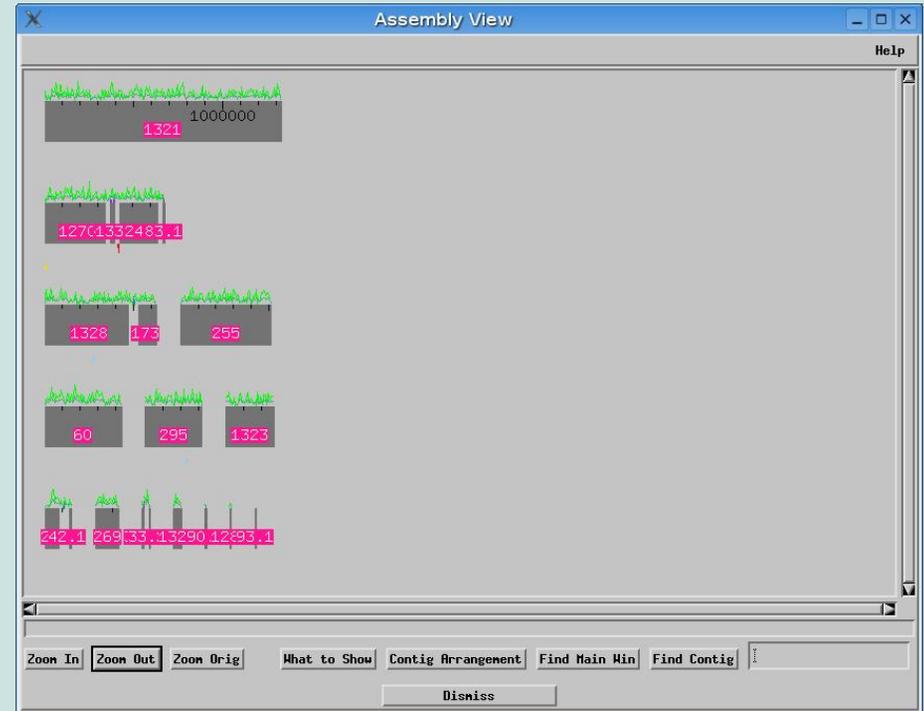
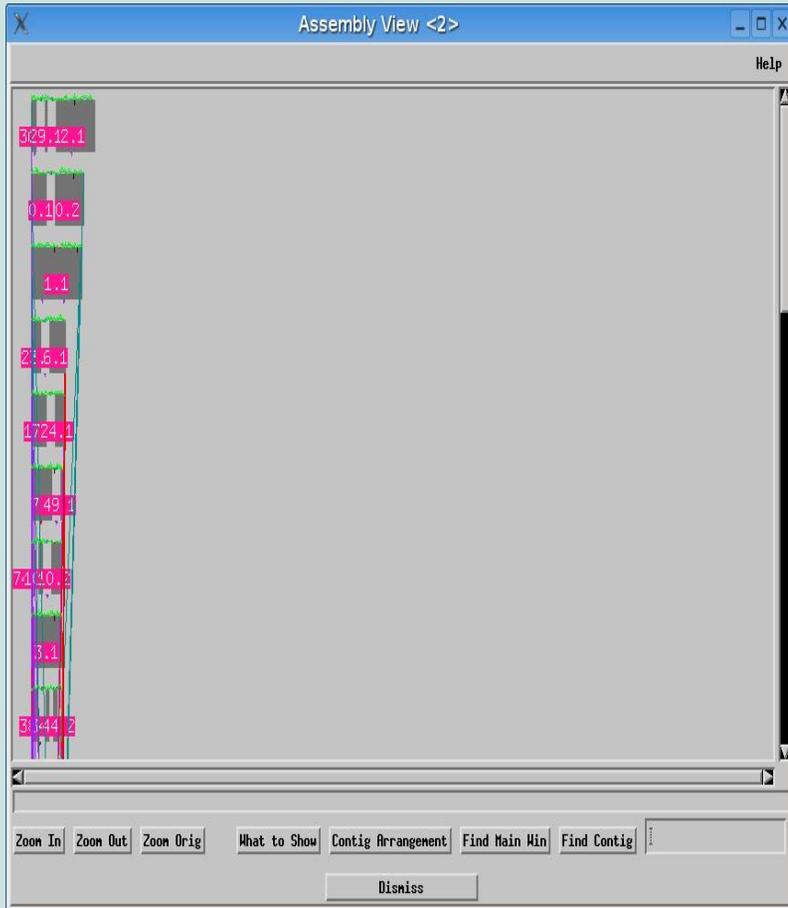
awollam@watson.wustl.edu

Results

- Thus far, twelve bacterial genomes have been improved
- Results show that MSI assemblies have longer N50 contig lengths, more contiguity, and more consistent read pair placement indicating better representation of the genome.

awollam@watson.wustl.edu

Increased contiguity (example1: BC genome)



MSI assembly: more contiguous

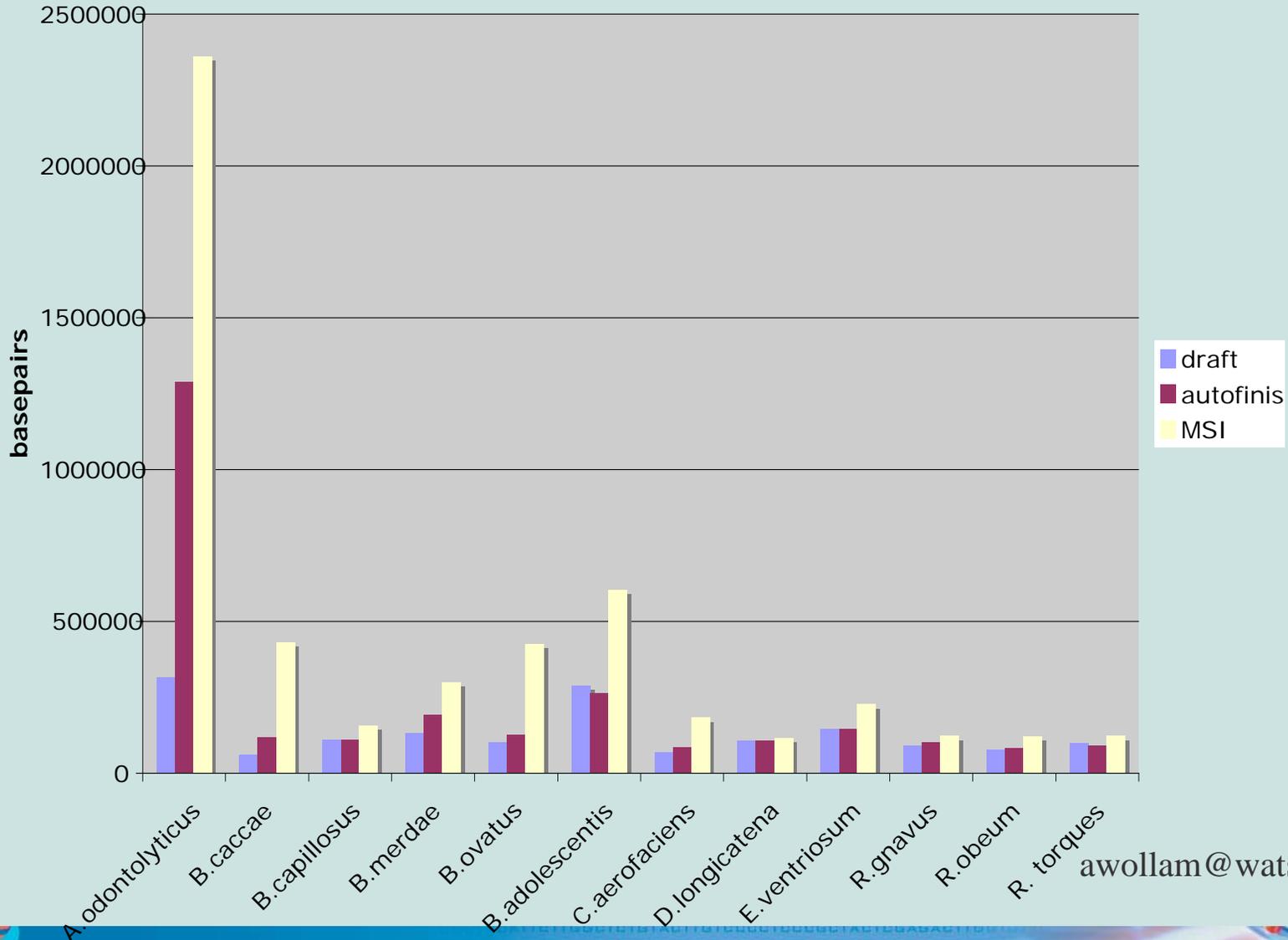
max contig: 1,329,168 bp

draft assembly: fragmented
maximum contig: 231,267 bp

awollam@watson.wustl.edu

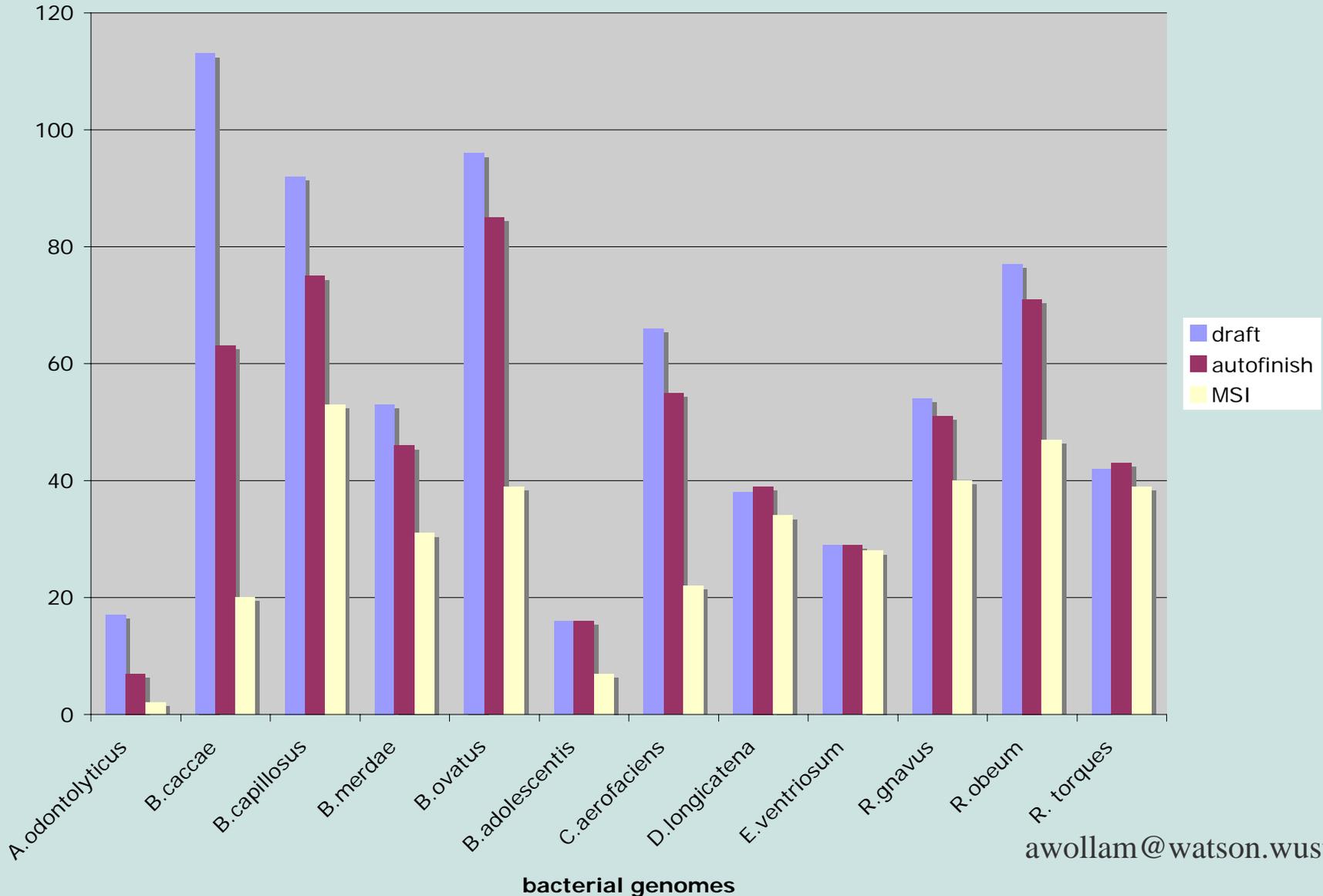
Comparison of assemblies at varying stages of improvement

comparison of N50 contig length



awollam@watson.wustl.edu

Comparison of assemblies at varying stages of improvement (total major contigs)



awollam@watson.wustl.edu

Summary of comparison between draft and MSI assemblies

- N50 contig length in twelve genomes = 1.2 - 7.5 folds increase
- average N50 contig length increase for twelve genomes = 3 folds
- Percentage decrease in total major contigs (more than 2 kb in size) = 3.5 -88.2%
- average decrease in contigs for twelve genomes = 44 %

awollam@watson.wustl.edu

finishing process versus whole genome MSI process

- more time needed as hard to finish regions can linger depending on the difficulty of genomes (5 Mb bacterial genome may take ~6 plus months)
- less high throughput (capacity of 2 Mb/month per finisher in BAC finishing)
- cost includes reagents, special techniques and more human resources.
- requires less time (5 Mb bacterial genome takes ~10 working days). Difficult regions are sorted as much as possible but not completely resolved.
- more high throughput (capacity of 20-25 Mb/month per finisher)
- approximately less than 10% of finishing cost (no additional reagent cost, less personnel)

awollam@watson.wustl.edu

annotation notes

- We will be looking at number of gene discovery, and annotation yield (confidence in functional assignments; number of functionally annotated versus conserved hypothetical and hypothetical proteins) as part of the parameters to be used to judge advantages of MSI.
- preliminary data on BC genome indicates MSI process eliminates false/split gene hits compared to draft assembly by 81.5%. (i.e. split genes were replaced by whole functional genes)
- 61% of draft genes are different than MSI genes, most of which go from smaller to larger sizes.
- detailed analysis of annotated regions between draft versus MSI assemblies are underway. Phylogenomic analysis of 1000 ORFS are currently in process.

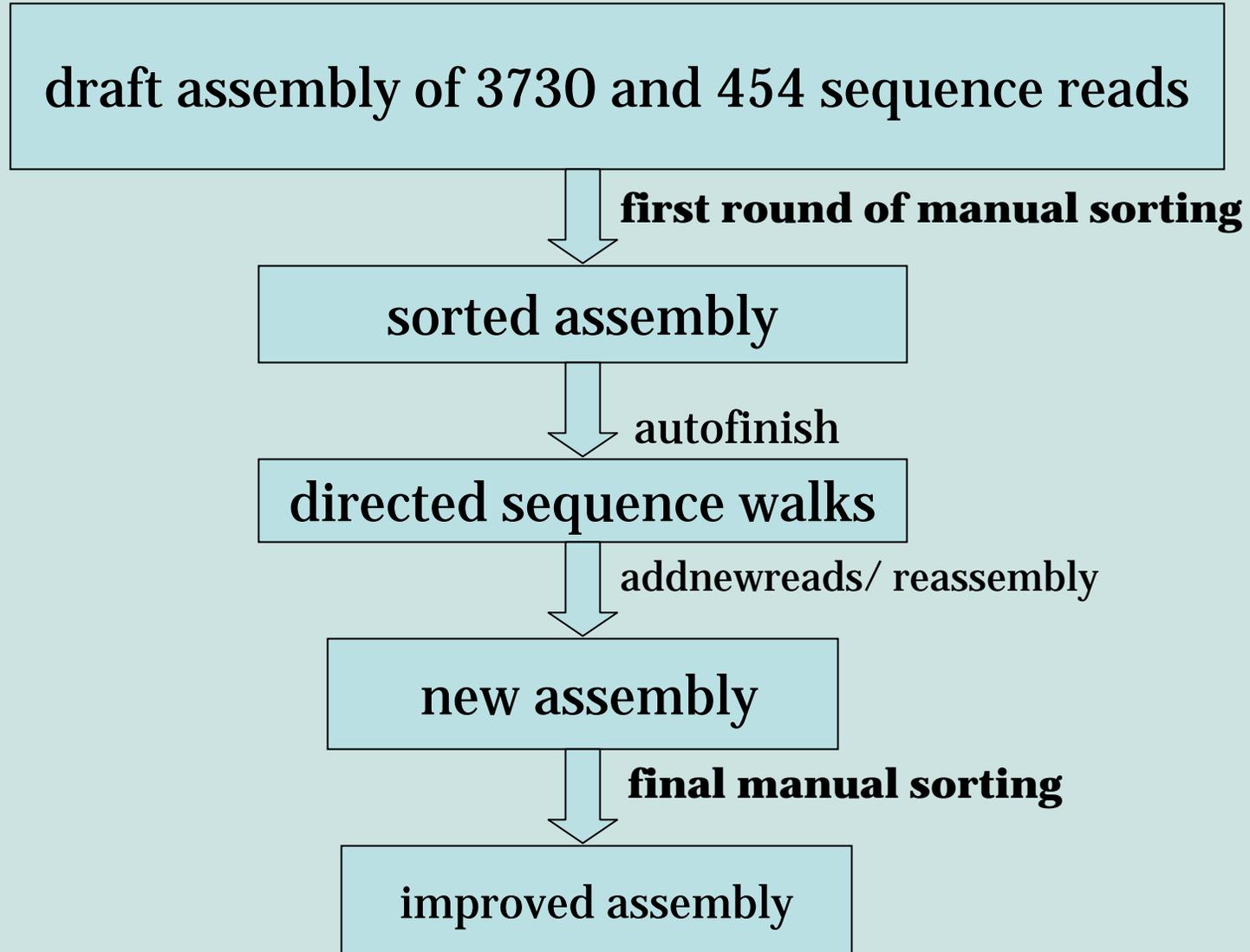
awollam@watson.wustl.edu

What's next?

- Modify MSI pipeline.
- Apply MSI process to large genomes for improving regions of interest.
- Auto-joining tools: development of software for determining possible joins between contigs, and completing joins.

awollam@watson.wustl.edu

modified MSI pipeline



summary

- Manual sequence improvement approach offers a better alternative to draft assemblies.
- To date, twelve bacterial genomes have been manually improved in significantly less time, cost and effort than those of finishing process.
- The result shows a sharp increase in N50 contig length, increased contiguity, less misassembled regions and better annotation when compared to the draft assembly.
- It is a better suited approach when dealing with large number of genomes

awollam@watson.wustl.edu

Whole genome finishing and manual improvement team

Bob Fulton

Aye Wollam

Neha Shah

Tom Wang

Matt Cordes

Kelsi Rotter (part time)

Jennifer Hodges (part time)

awollam@watson.wustl.edu

