

The “Art” of Shotgun Assembly

PacBio Bioinformatics Developer Meeting, Aug 26, 2015

FIND MEANING IN COMPLEXITY

Large Genome Assembly: A Bit of Every Everything

Information Technology

- Grid computing
- Cloud computing
- Distributed file system
- Data management

Software Engineering

- Open source development
- Modularization
- Workflow management
- Deployment

easy

Sequence Data Analytics

- Data quality assessment
- Genome repeat content estimation
- Algorithm parameter optimization
- Effective data reduction

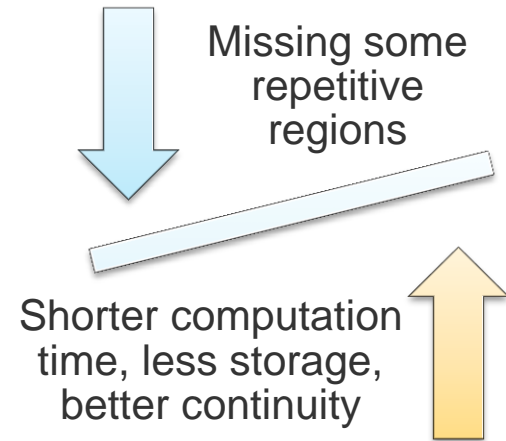
Genomics Algorithms

- Overlappers
- Assembly graph construction and reduction
- Metadata tracking and management
- Understand the trade-off between various factors

Current Challenges And Solutions (for Falcon)

- Reduce false overlapping caused by repeats
 - What is the right amount of data for XYZ?
 - Read length cutoff?
 - Sensitivity and specificity?
 - One month or 48 hours computation?
 - Find and understand “good rules” for assembly graph reduction
 - Overlap filtering
 - Remove “repeat induced bridging unitigs”
-
- Diploid genome “unzipping”
 - New concept to adapt (<https://speakerdeck.com/jch1n>)
 - Need to track every raw read
 - “Augmented alignment” needed for Quiver polishing

Main challenge: balance the trade-off



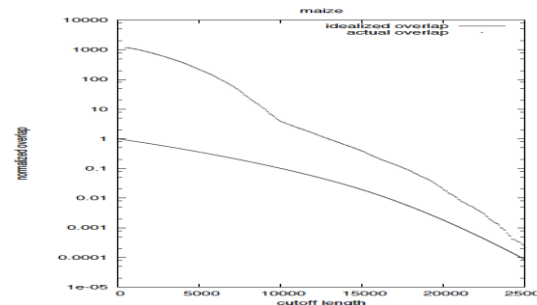
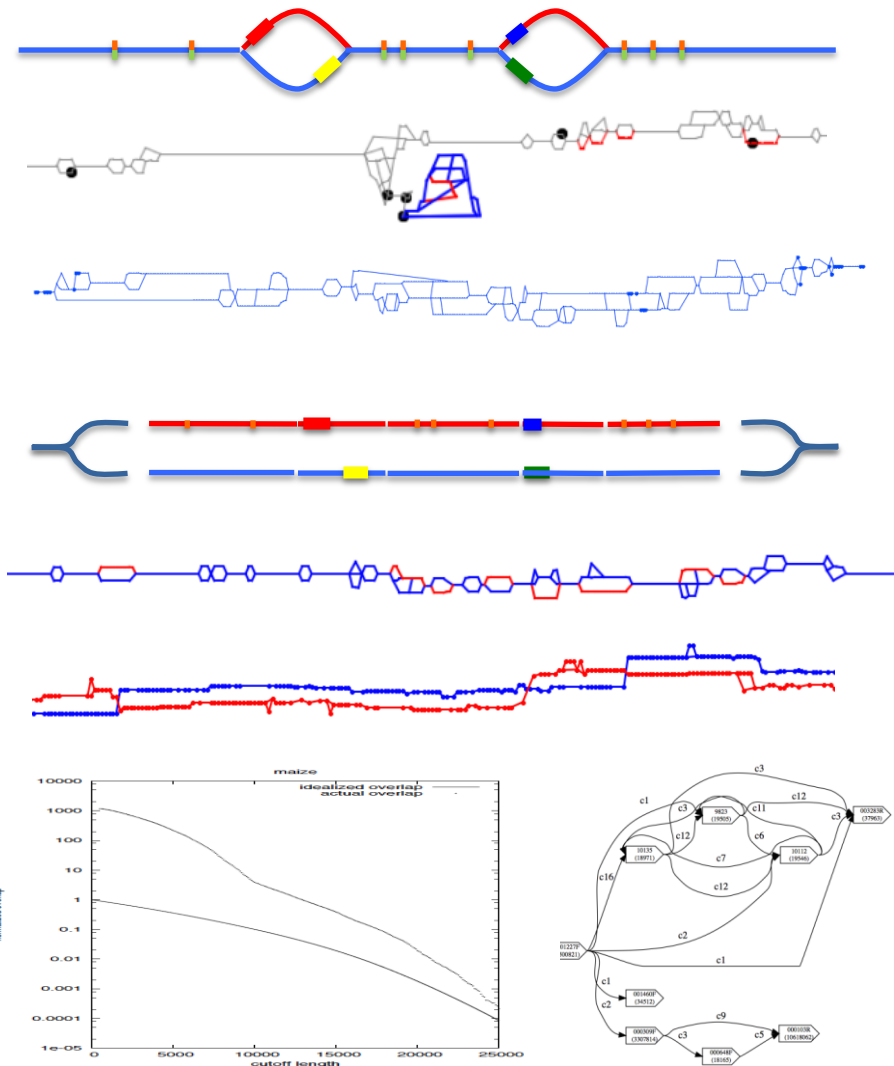
A better understanding of genome **repeat structures** is the key!!

Challenges

- More complicated data model
- Evaluation and validation can be tricky

Near Future Opportunities

- Assembly graph visualization GUI for manual refinement or assembly error detection
- Estimate optimized assembly algorithm parameters by analyzing small sample for “**REPEATREPEAT**”
- Multi-step assembly to provide full pictures including all repeat levels
- Algorithm and efficiency improvement for polyploid genomes
- Combination of other types of data to improve the final assembly quality





PACIFIC
BIOSCIENCES®

For Research Use Only. Not for use in diagnostic procedures. Pacific Biosciences, the Pacific Biosciences logo, PacBio, SMRT, SMRTbell, and Iso-Seq are trademarks of Pacific Biosciences. BluePippin and SageELF are trademarks of Sage Science. NGS-go and NGSengine are trademarks of GenDx. All other trademarks are the sole property of their respective owners.