LA-UR-12-10336

Title: Genome Improvement with PacBio sequencing technology

Author(s): Han, Shunsheng
Chertkov, Olga
Daligault, Hajnalka E.
Davenport, Karen Walston
Detter, John C.
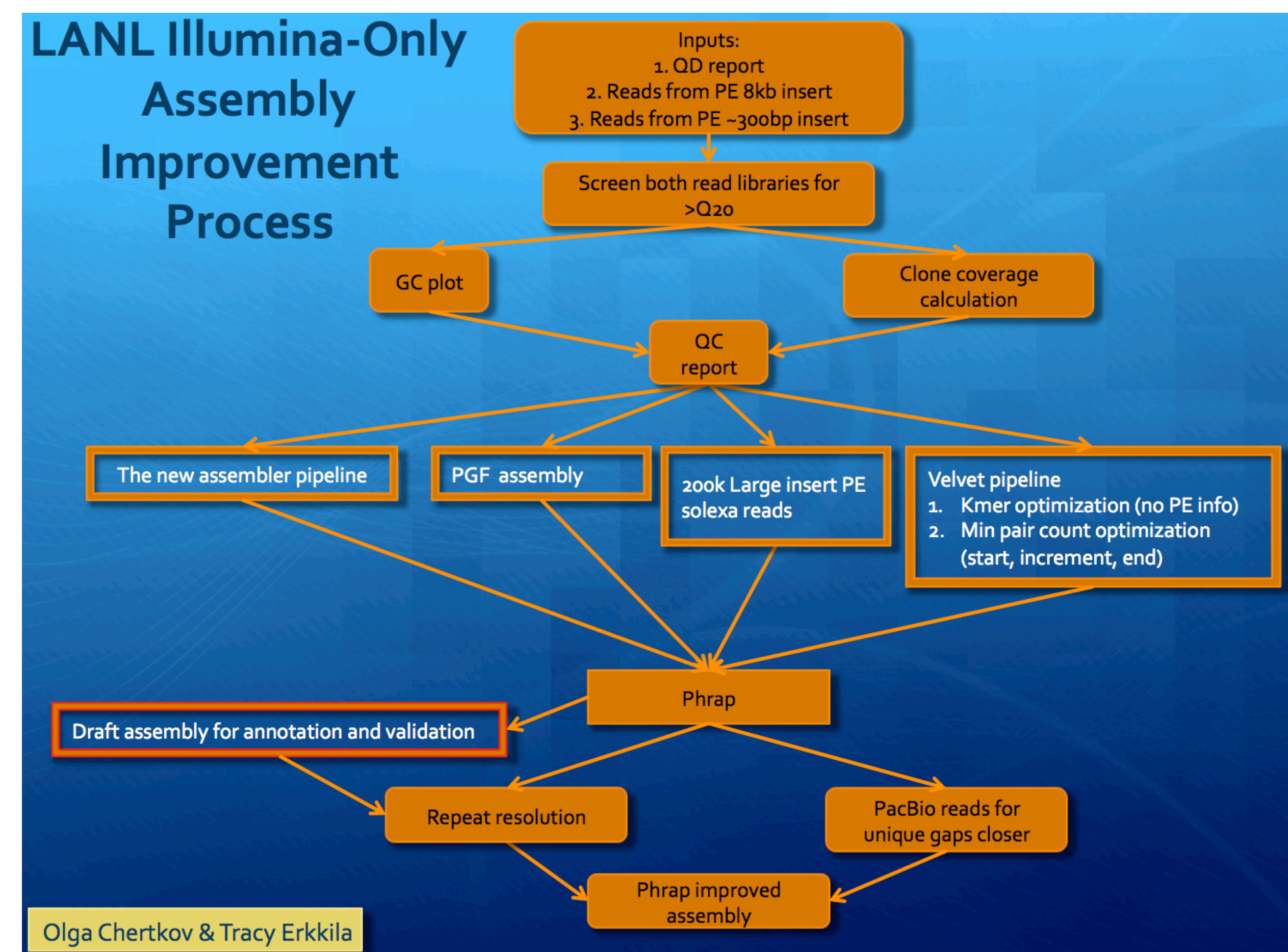Munk, A. Christine
Reitenga, Krista G.
Zhang, Xiaojing

Intended for: DOE
AGBT, 2012-02-15/2012-02-18 (Marco Island, Florida, United States)
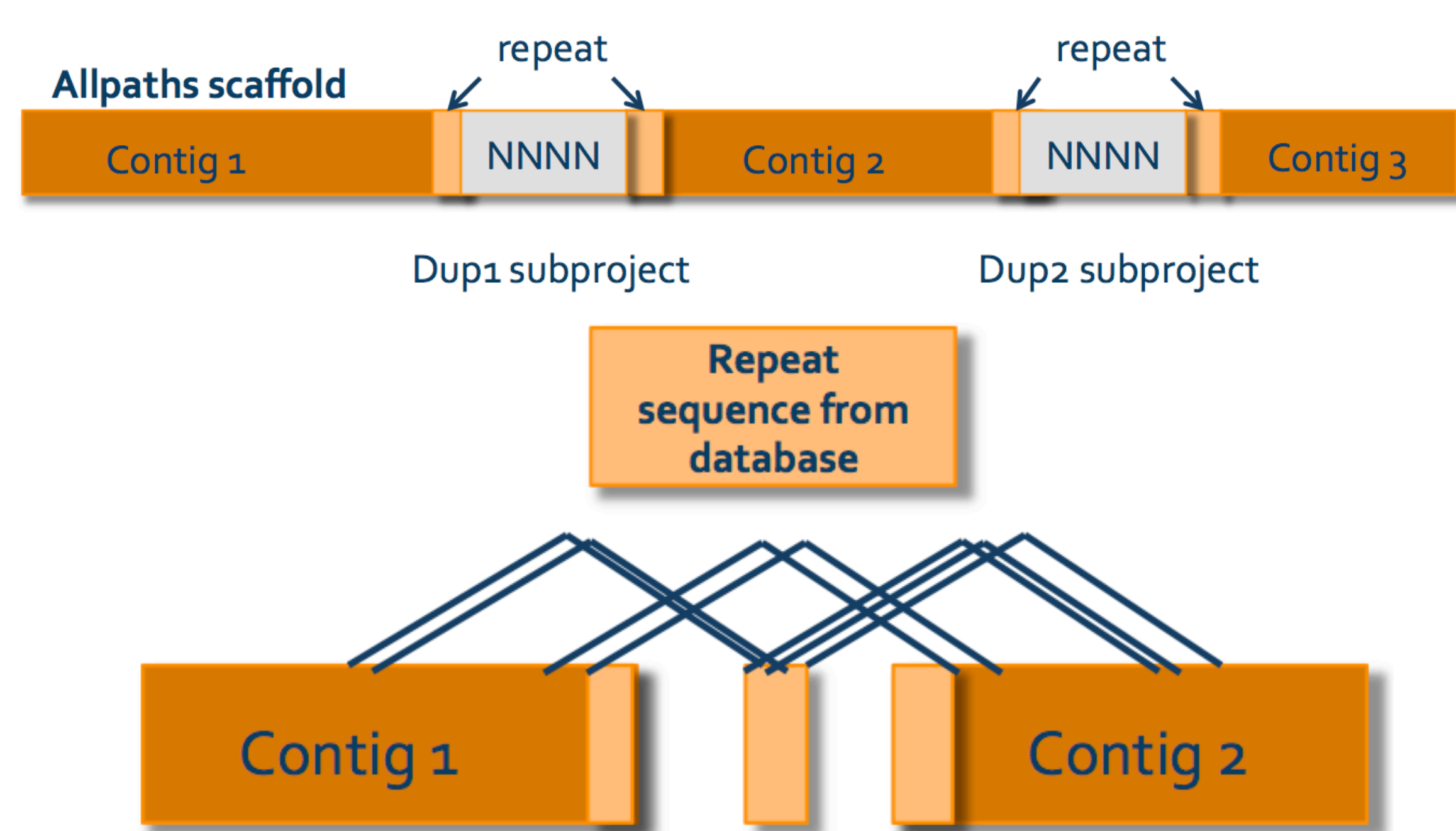NMED
US EPA
Biological resources
Reading Room
RCRA

![Los Alamos National Laboratory logo] Los Alamos
NATIONAL LABORATORY
—— EST. 1943 ——

# Genome Improvement with PacBio sequencing technology

Cliff S. Han[1], Lucy (Xiaojing) Zhang[1], Krista Reitenga[1], Primo Baybayan[2], Susana Wang[2],
Karen W. Davenport[1], Hajnalka E. Daligault[1], Olga Chertkov[1], Chris Munk[1], and Chris Detter[1]

1 The DOE Joint Genome Institute, Los Alamos National Laboratory, Los Alamos, NM 87545; 2 Pacific Bioscience, 1380 willow Rd. Menlo park, CA 94025.

LA-UR 12-01419

**Abstract**: As massively parallelized sequencing technologies enable draft genome sequencing at minimal cost, the cost for closing gaps in the draft genome with current Sanger methods becomes disproportionally high. The DOE Joint Genome Institute sequences hundreds of bacterial genomes each year. There could be dozens to hundreds of contigs in draft assemblies produced with only Illumina data, which at JGI include short (~300 bases) and medium (~8 Kb) insert libraries. A significant of them need to be finished/improved to a better level of continuity. Data from PacBio has been proven to improve continuity of assembly as it has better coverage in high or low GC regions comparing to Illumina technology. Unfortunately, the throughput of the RS currently does not compare to that of our Illumina pipeline yet. An alternate approach is to sequence the gaps and repetitive regions of genomes drafted with Illumina technology. Amplicon sequencing with PacBio RS machine has the potential to replace our current Sanger based genome finishing process. We used two finished genomes with high and medium GC content as test cases. Illumina data was reproduced with current technology with libraries of short and medium insert sizes. The two assemblies of the Illumina-only data produced 98 and 105 contigs in a limited number of scaffolds. PCR products from those gaps were pooled and sequenced with a PacBio machine. Subreads were aligned to known sequences next to PCR primer regions. Subreads that belong to a single PCR amplicon were assembled separately. The assembled sequences were aligned with the primer sequences again. Sequences between the primer pair are used to close gaps in the main project. Hard stops with 30 bp hairpin structure can be sequenced with PacBio without problem.
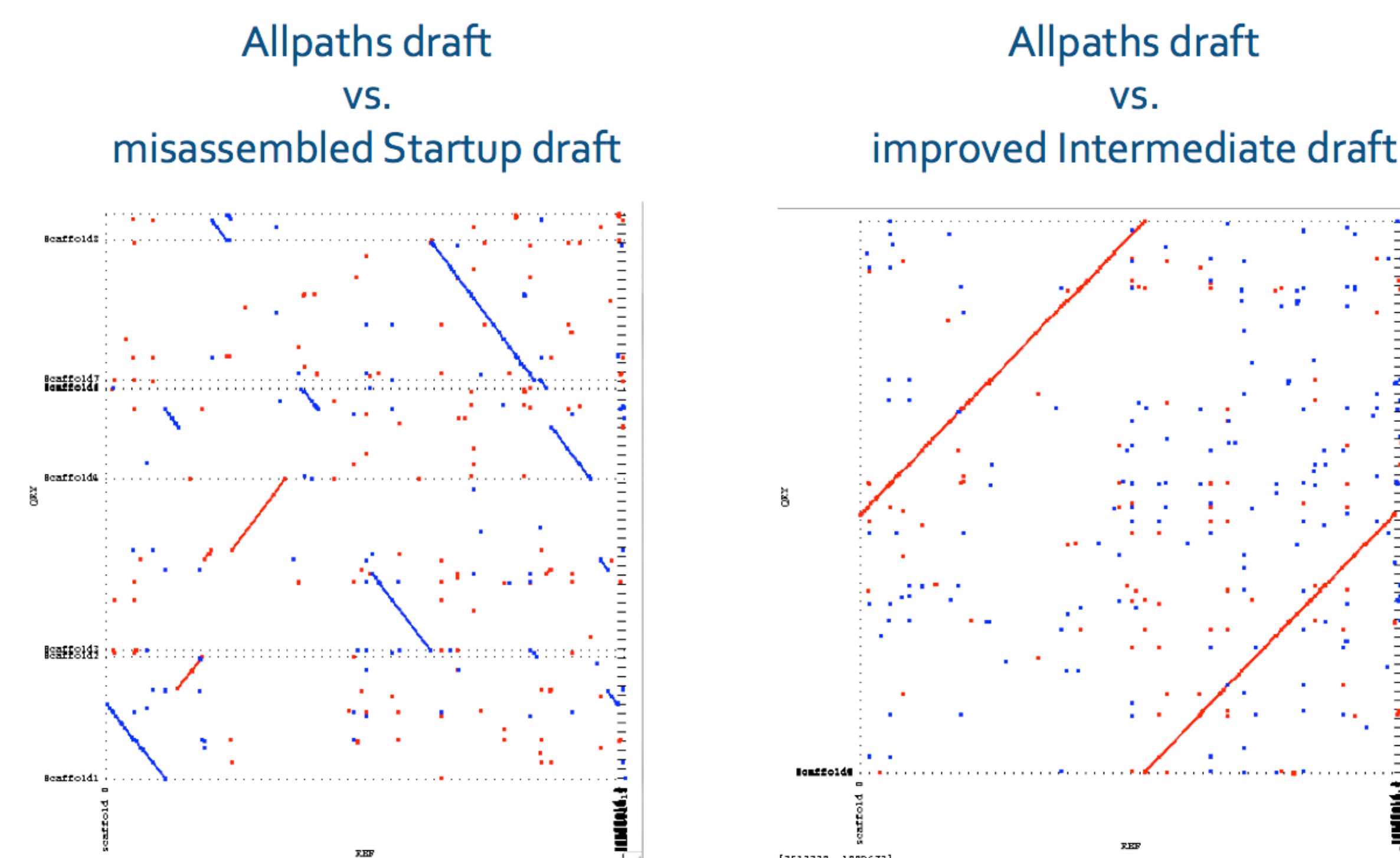
## Repeat Resolution software
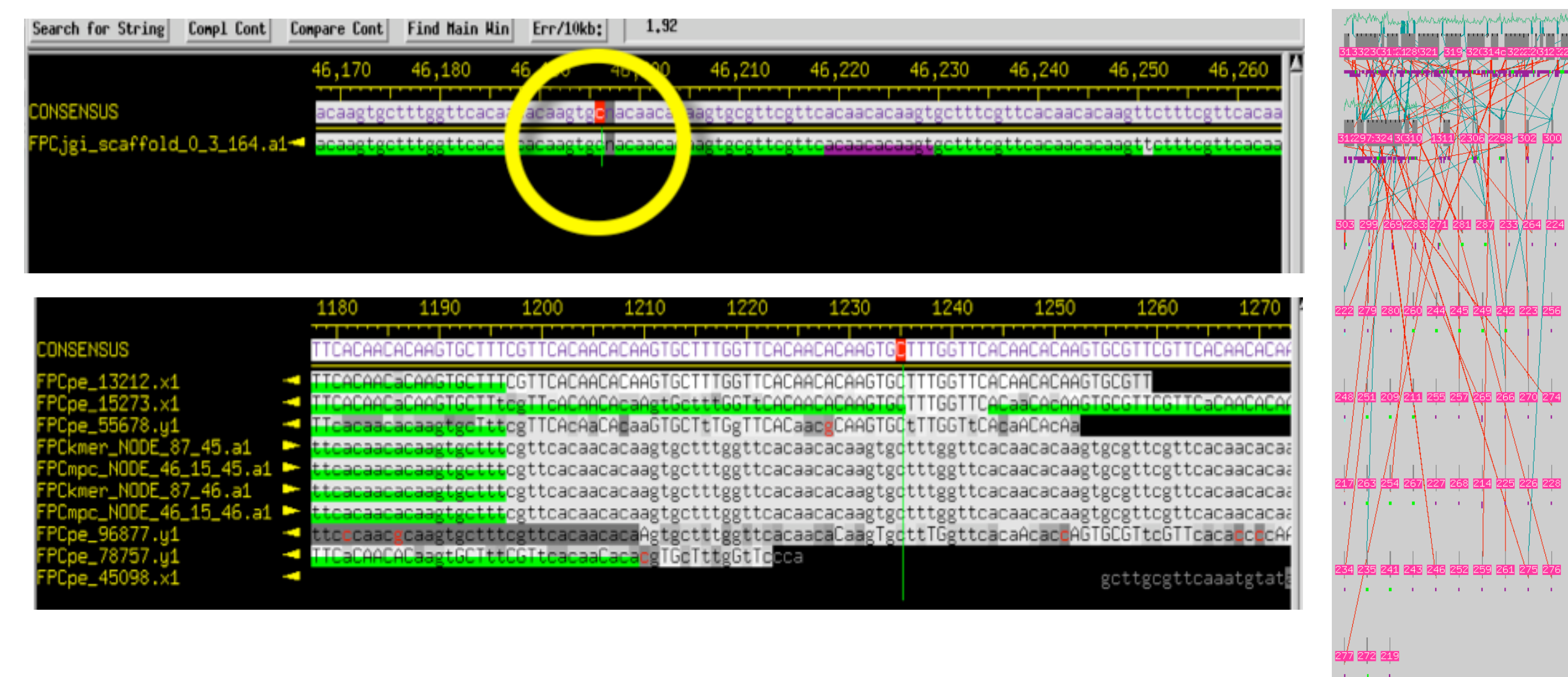


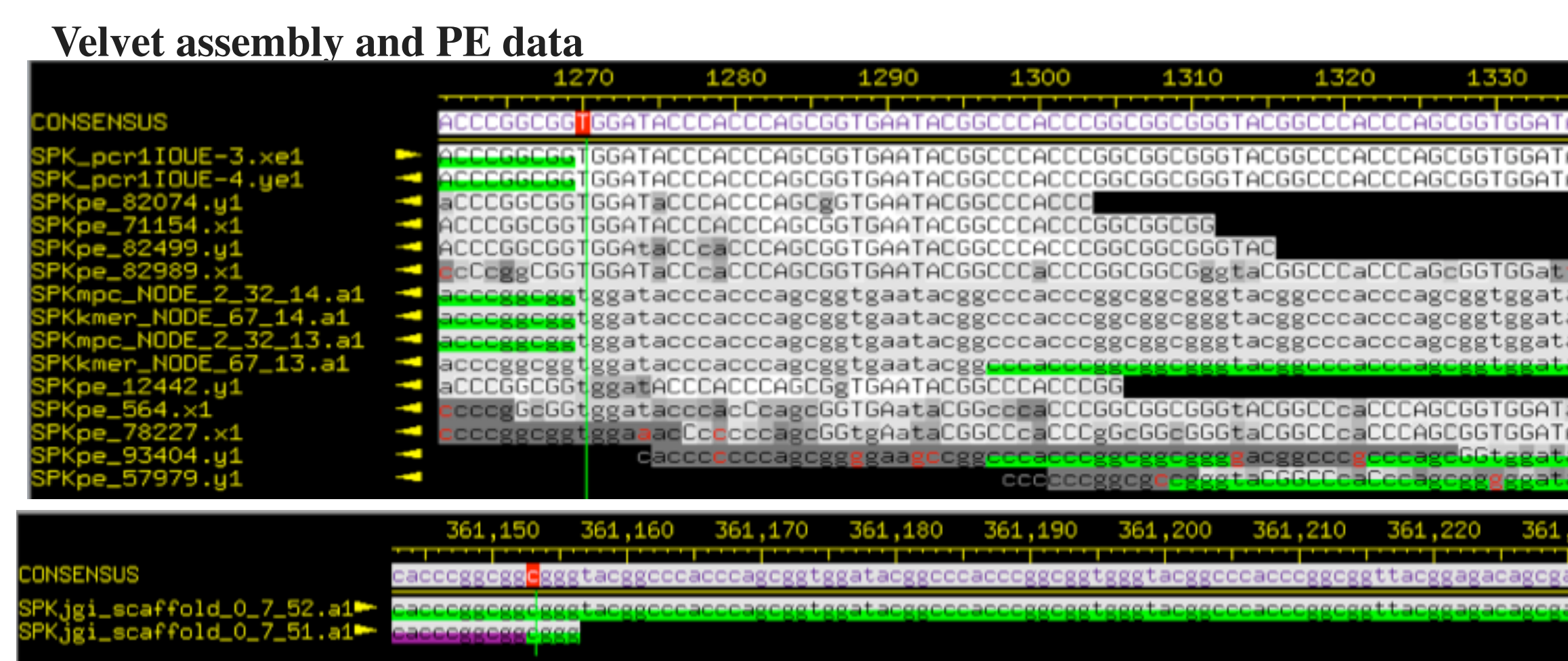## Scaffolds: Allpaths is more accurate

Allpaths draft vs. misassembled Startup draft

Allpaths draft vs. improved Intermediate draft



## Sequence: Velvet is more accurate



Missing sequence in Allpaths assembly

Velvet assembly and PE data

Allpaths- missing sequence

## LANL Illumina-Only Assembly Improvement Process



Olga Chertkov & Tracy Erkkila

## Genome improving pipeline
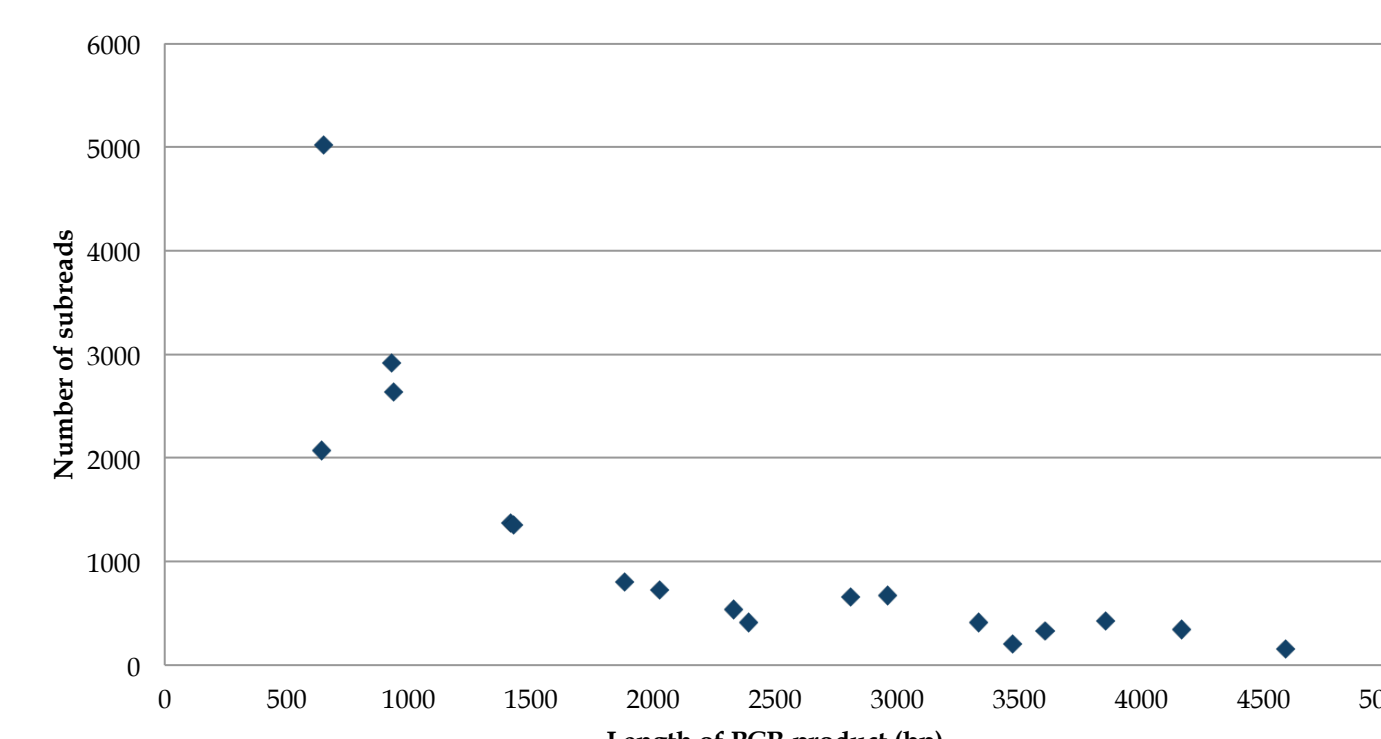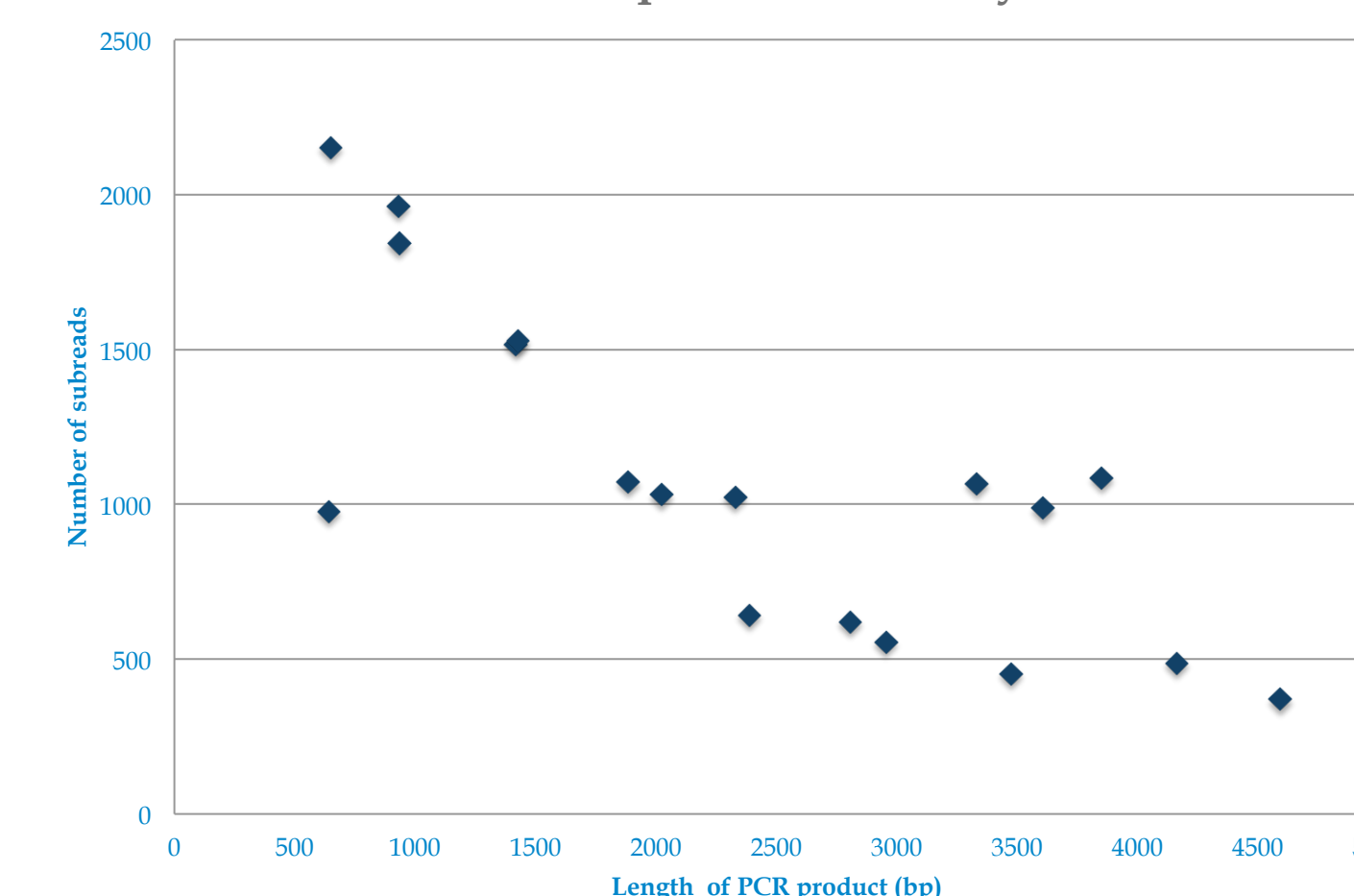


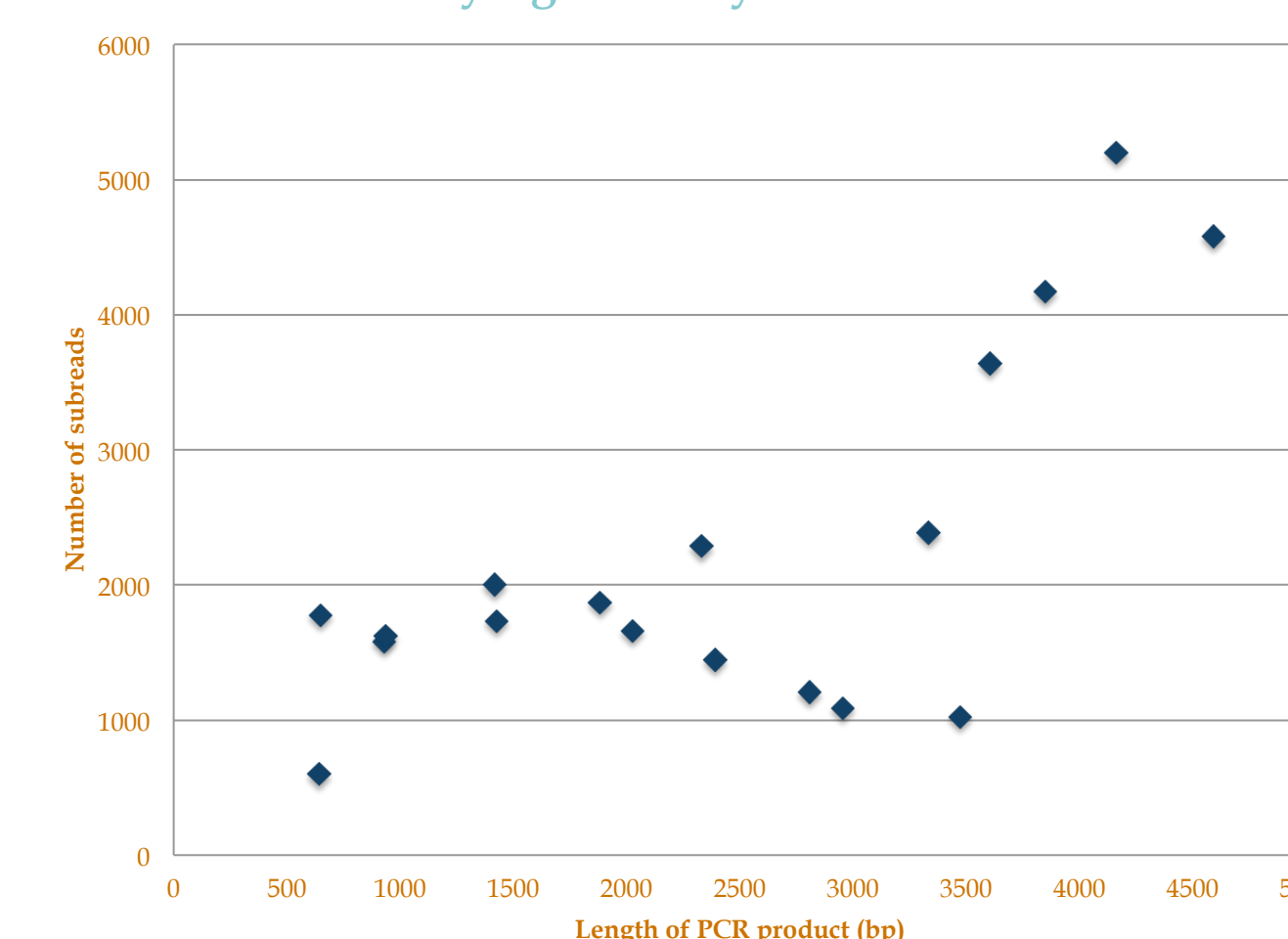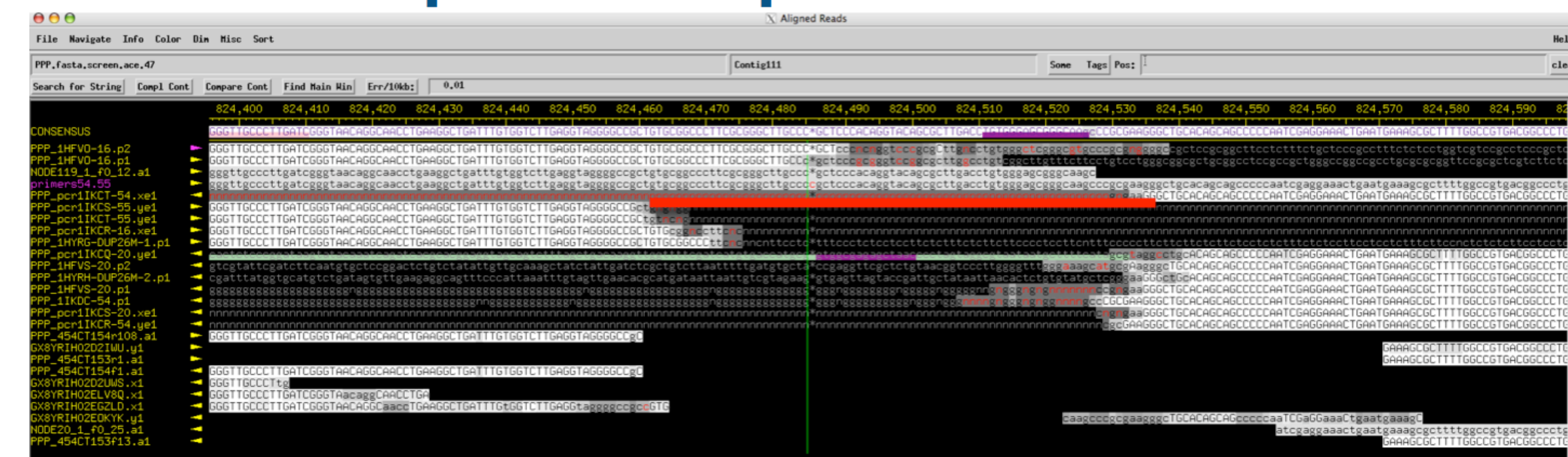## PacBio Sequencing of PCR for Gap Closure



Trial 1: Equivalent Mass

Trial 2: Equivalent Molarity

Trial 3: Varying Molarity Based on PCR size

## Hard stops are no problem with PacBio



cccttcgccgggcttgcccgctcccacaggtacagccgcttgacctgtgggagcgggcaagcccgcgaaggg

30bp long | 30bp long

hairpin loop

*An example of a sequence that creates a hairpin loop in the secondary structure which has been very difficult to complete when sequencing a genome.*

Shown in the sequences above,
only the PacBio read (underlined in red) goes through the hard stop region.
Illumina, 454, and Sanger (PCR) reads didn't.