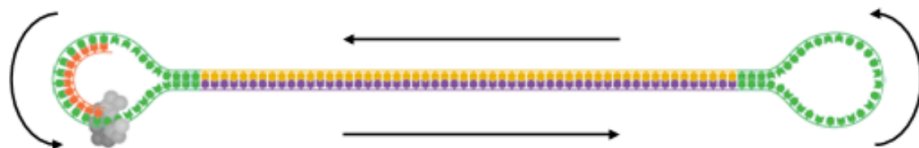


## Long read Assembly



- The main protocols for long read technologies
- The error profile of long reads
- Methods of read correction
  - Correction with short reads
  - Correction using long reads
- Assembly Tools
- Assembly Diagnostics
- Assembly Polishing



**SMRTbell™ Template**



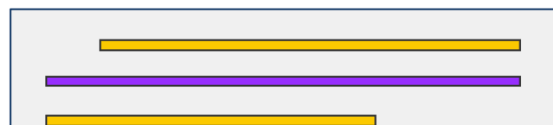
## Polymerase Read

### Definition:

- Sequence of nucleotides incorporated by polymerase while reading a template
- Includes adapters
- Often called “read”
- Includes adapters
- 1 molecule, 1 pol. read

### Purpose:

- QC of instrument run
- Benchmarking



## Subread

### Definition:

- Single pass of template
- Adapters removed
- 1 molecule,  $\geq 1$  subreads

### Unique data:

- Kinetic measurements
- Rich QVs

### Purpose:

- For subsequent analysis



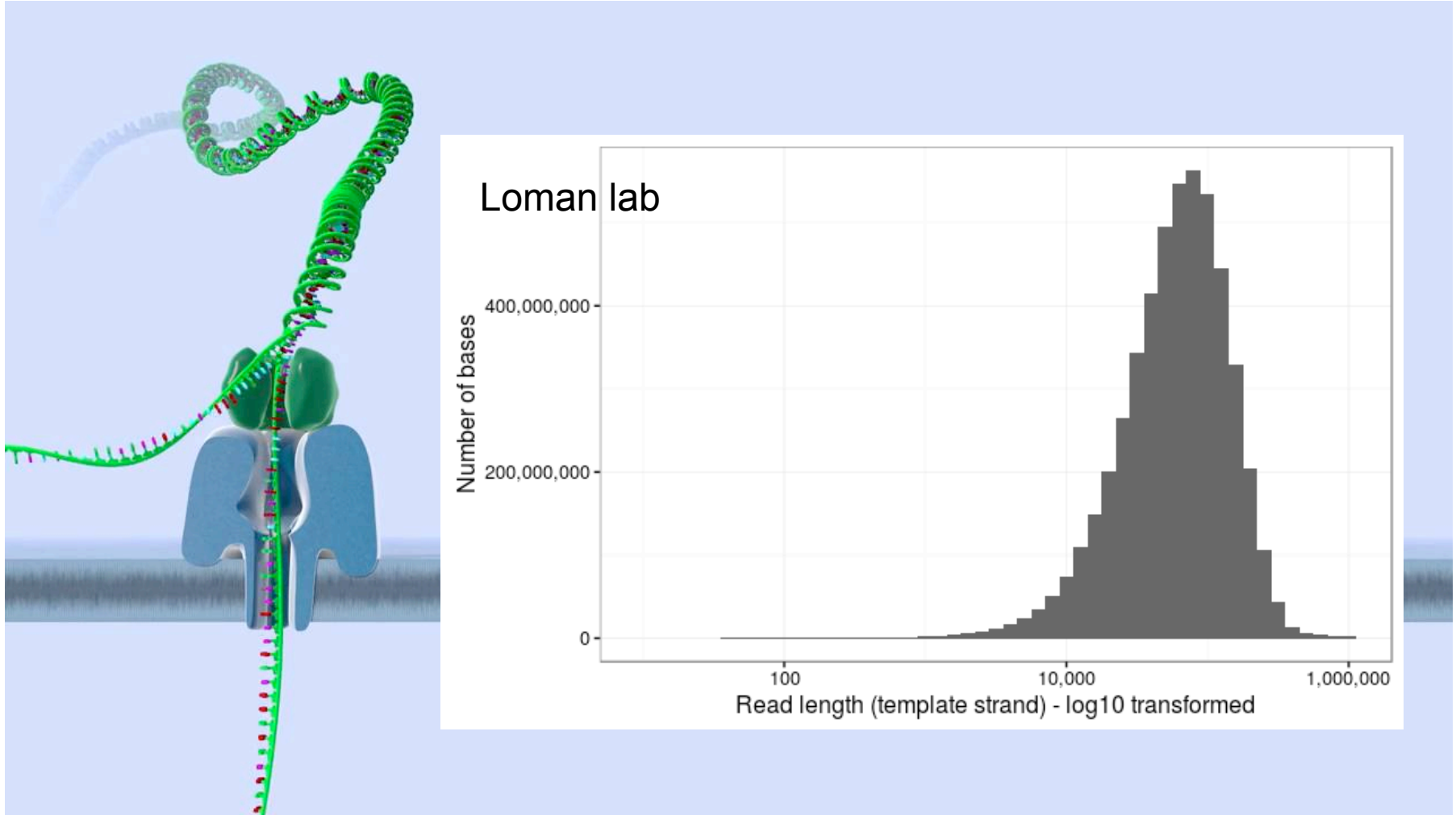
## Read of Insert

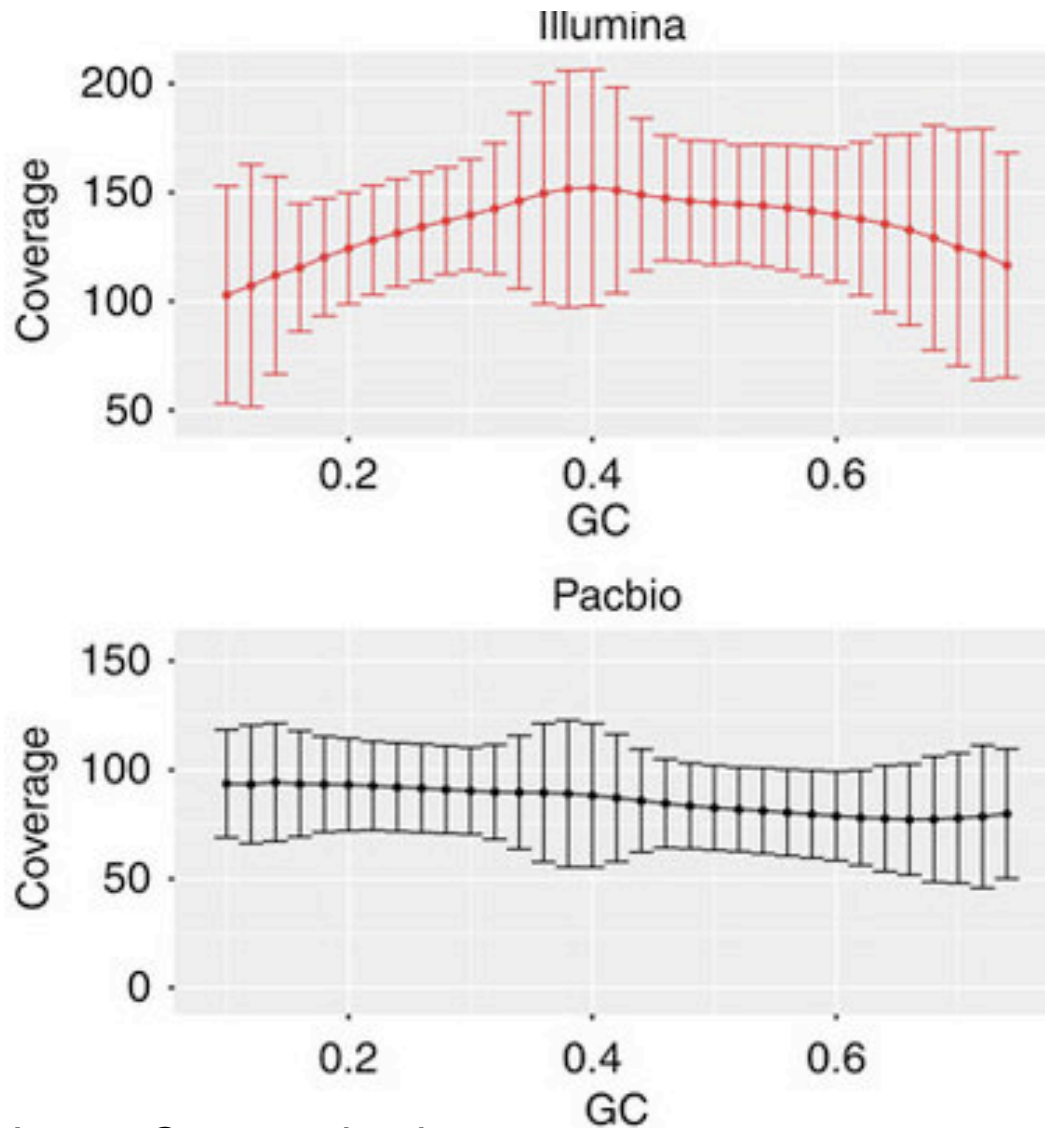
### Definition:

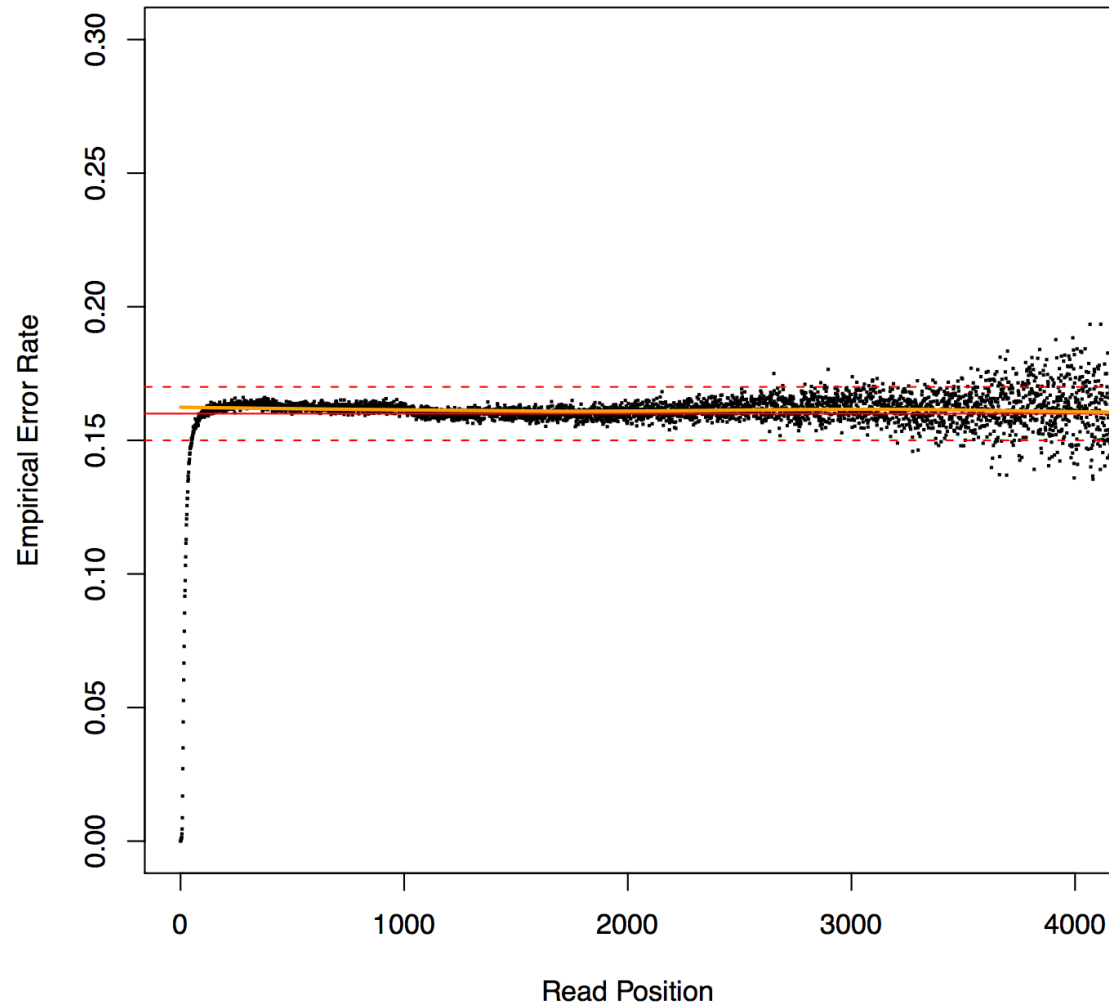
- Represents highest-quality single-sequence for an insert, regardless of number of passes
- Generalizes CCS for  $< 2$  passes and  $RQ < 0.9$
- 1 or more passes
- 1 molecule, 1 read

### Purpose:

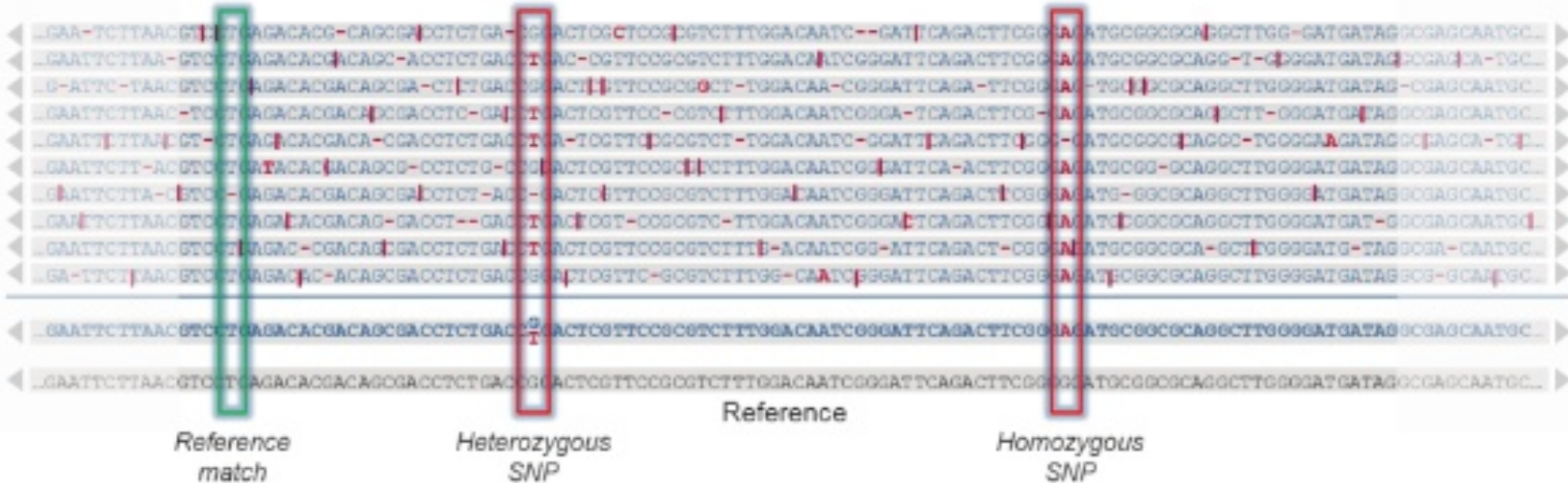
- For Library QC
- For subsequent analysis







## PacBio: error rate

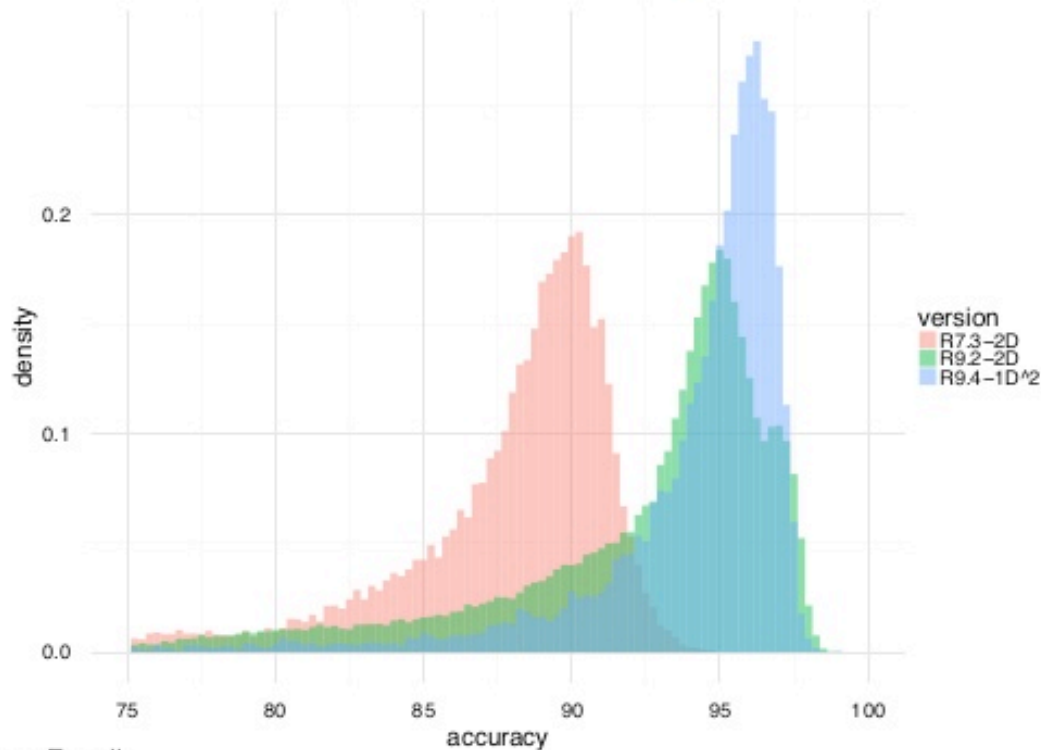


Single read: 86%

30x Consensus: 99.999%

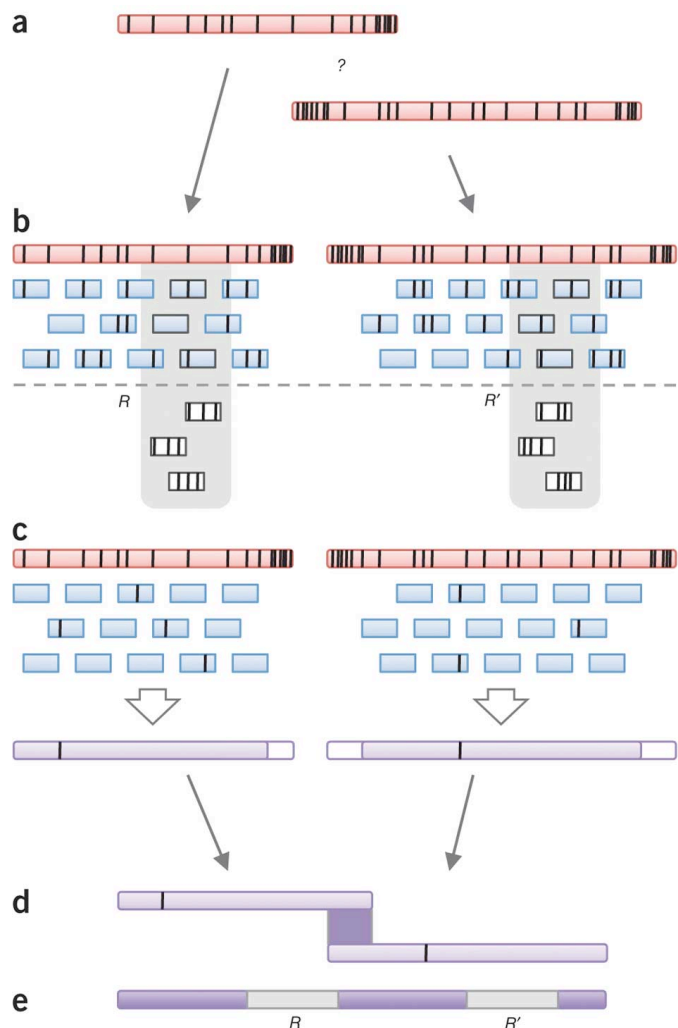


## 1D<sup>2</sup> Accuracy



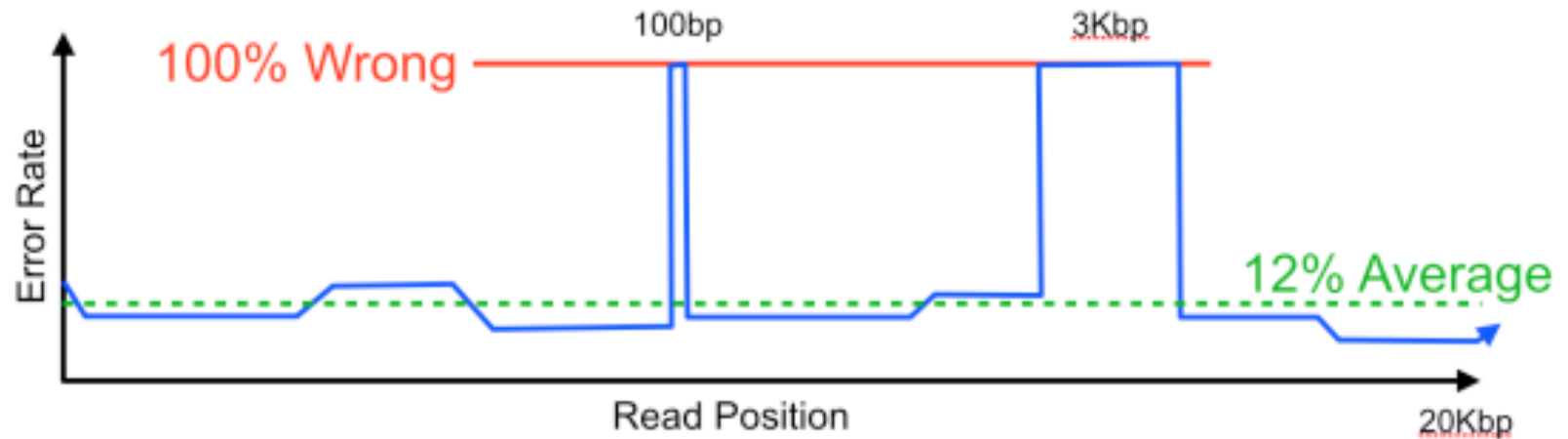
- All runs are E. coli
- R7.3-2D from Nick Loman
- R9.2-2D from OICR
- R9.4-1D<sup>2</sup> provided by ONT

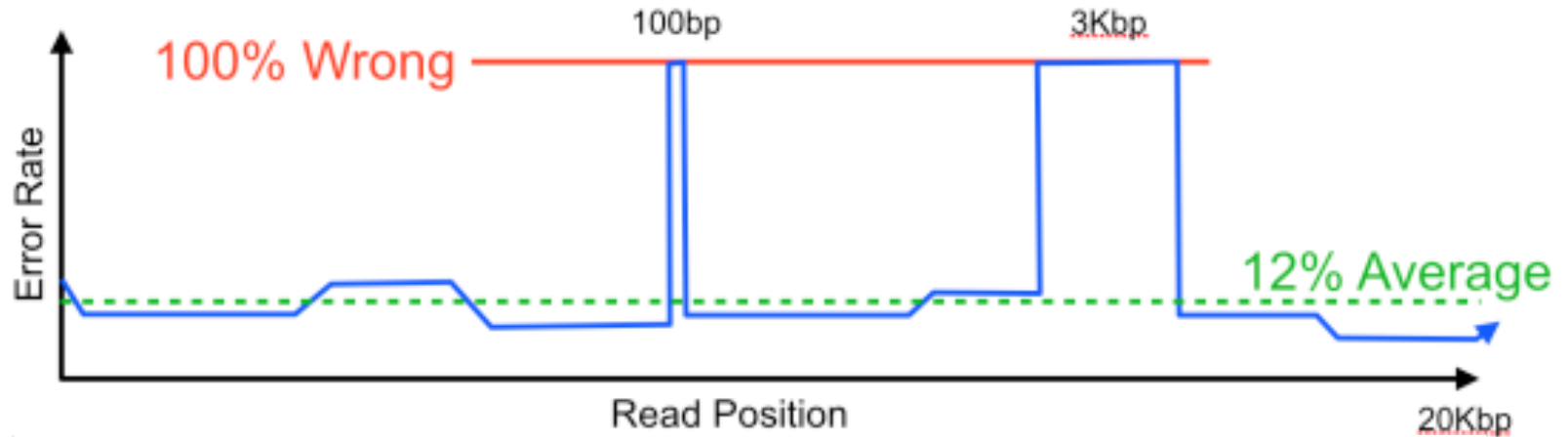




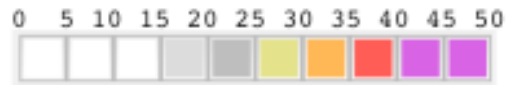
- Correction using Illumina reads
- Homopolymer correction, point mutations, and indels
- Doesn't correct structural errors
- Tools
  - PBcR / PacBioToCA
  - LSC / LSCplus
  - LoRDEC (de Bruijn graph)
  - Proovread
  - ECTools
  - Jabba (de Bruijn graph)

- Structural errors
  - Chimeric reads ( see Tallon et al. 2014. BMC Genomics)
  - Missed or incorrectly inferred adapter
  - Interference from other molecules
  - ...

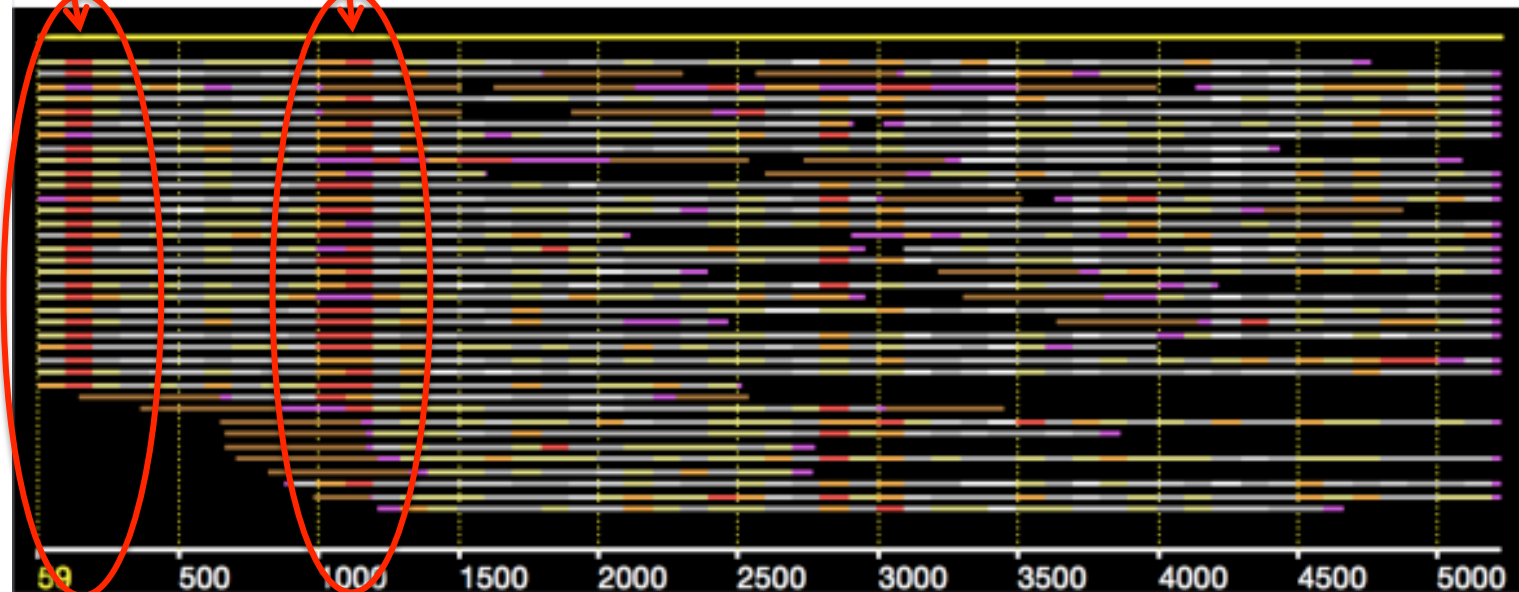




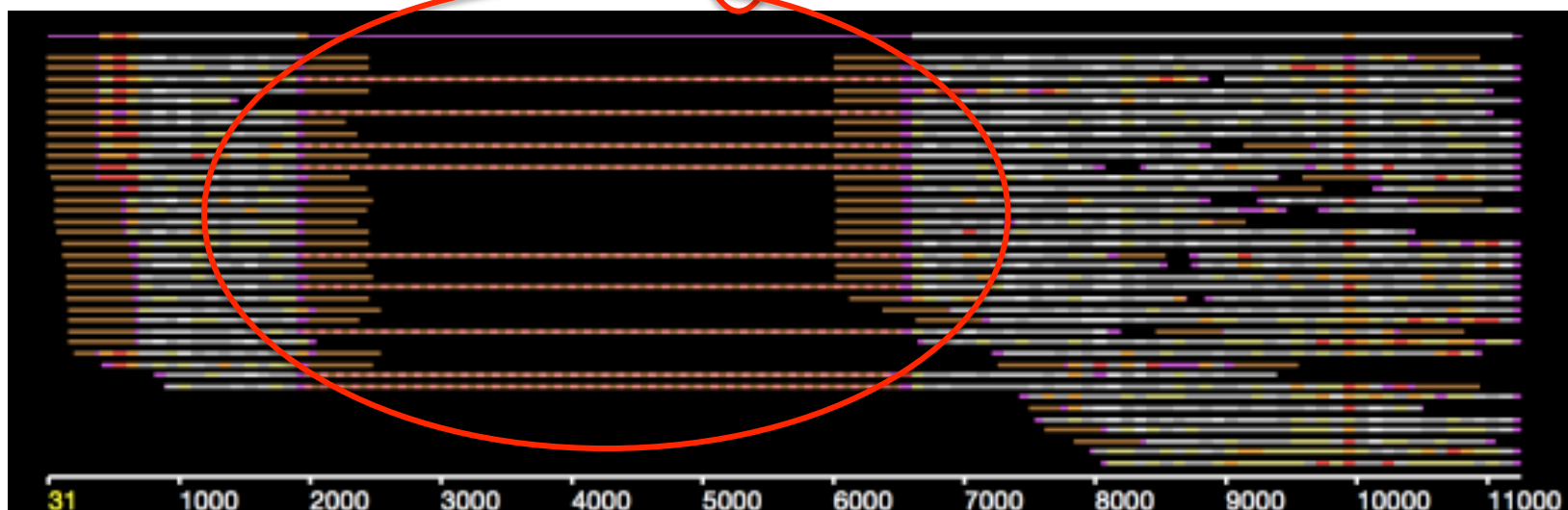
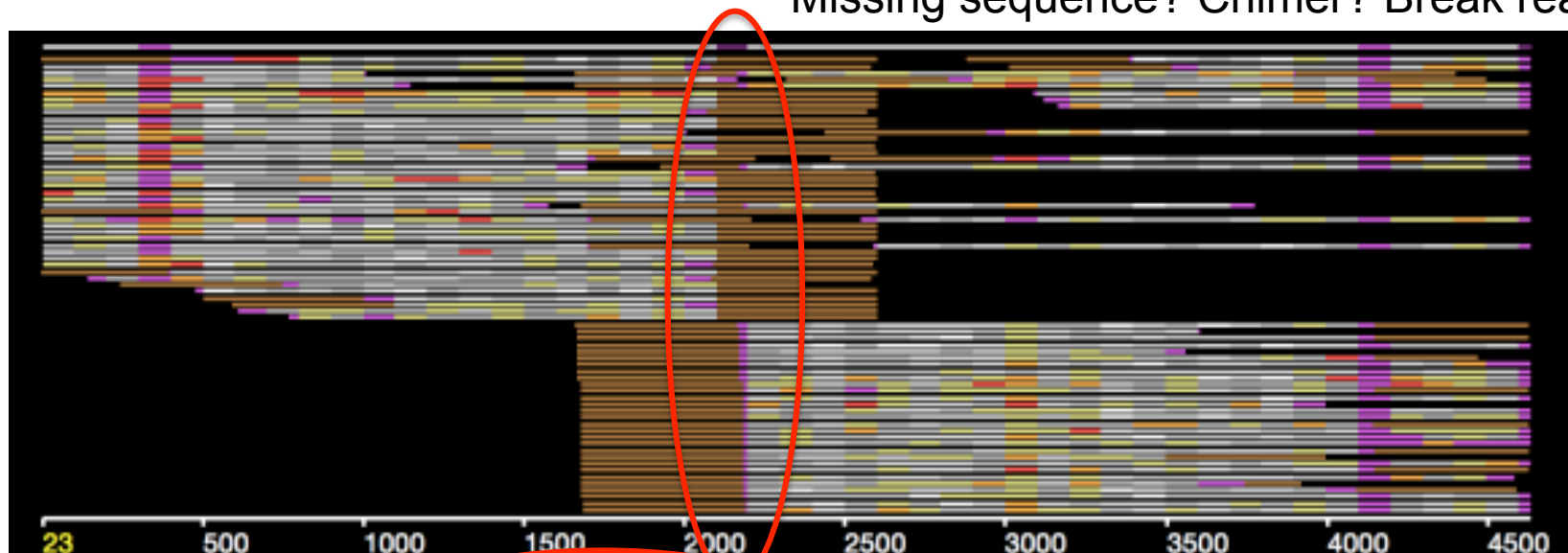
Likely wrong sequence



Differences per 100bp segment

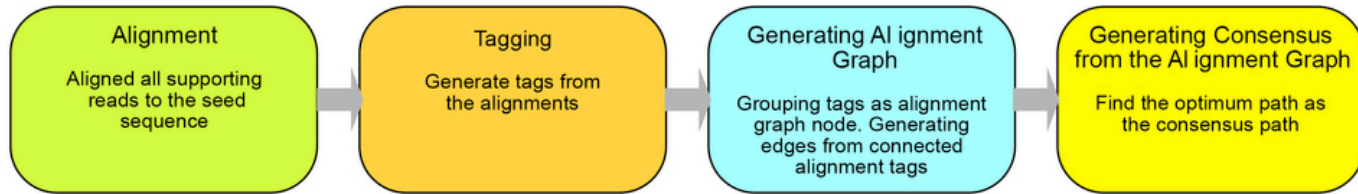


Missing sequence? Chimer? Break read



Inserted sequence from Chimer?

- Use a weighted directed acyclic graph to find consensus sequence



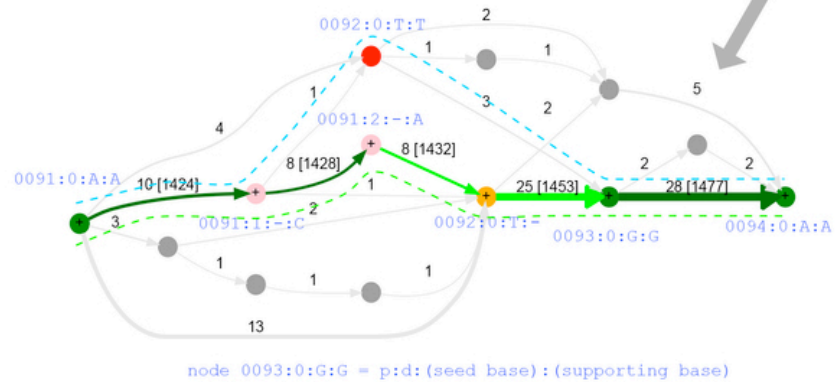
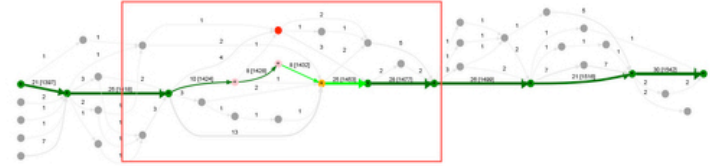
```

p =
d =
seed = ATATTA-GGC
read 1 = ATAT-ACGGC

p = 0123455678
d = 0000001000
seed = ATATTA-GGC
read 1 = ATAT-ACGGC

p =
d =
seed = AT-ATTA--GGC
read 2 = ATCAT--CCGGC

p = 011234555678
d = 001000012000
seed = AT-ATTA--GGC
read 2 = ATCAT--CCGGC
  
```



Seed Sequence Path	Consensus Sequence Path
0091:0:A:A	0091:0:A:A
0092:0:T:T	0091:1:-:C
0093:0:G:G	0091:2:-:A
0094:0:A:A	0092:0:T:-
	0093:0:G:G
	0094:0:A:A

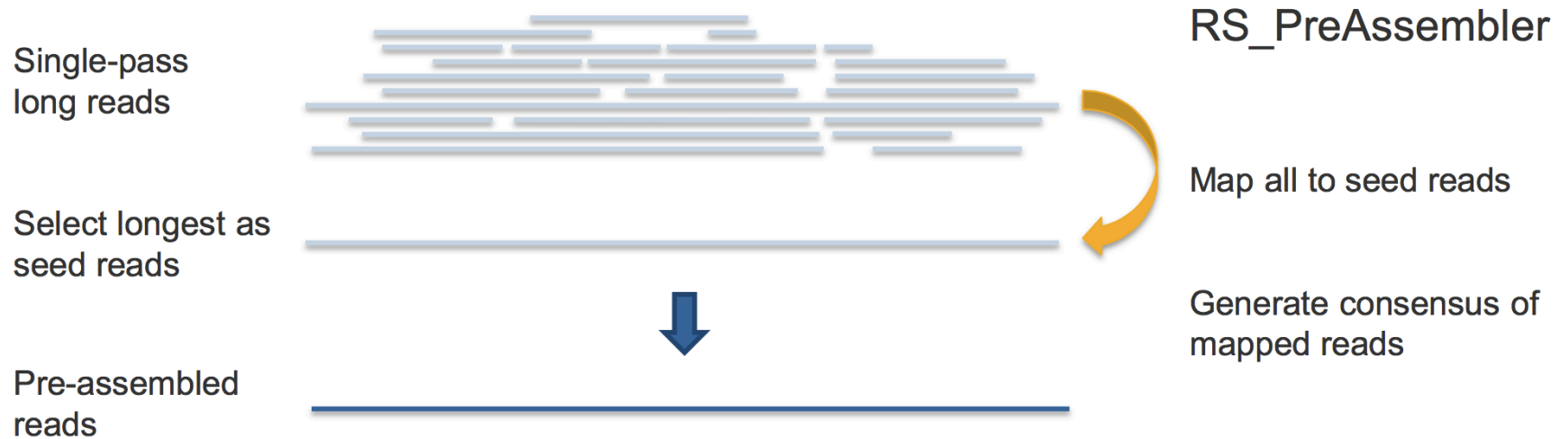
Seed Sequence	Consensus Sequence
ATGA	ACAGA

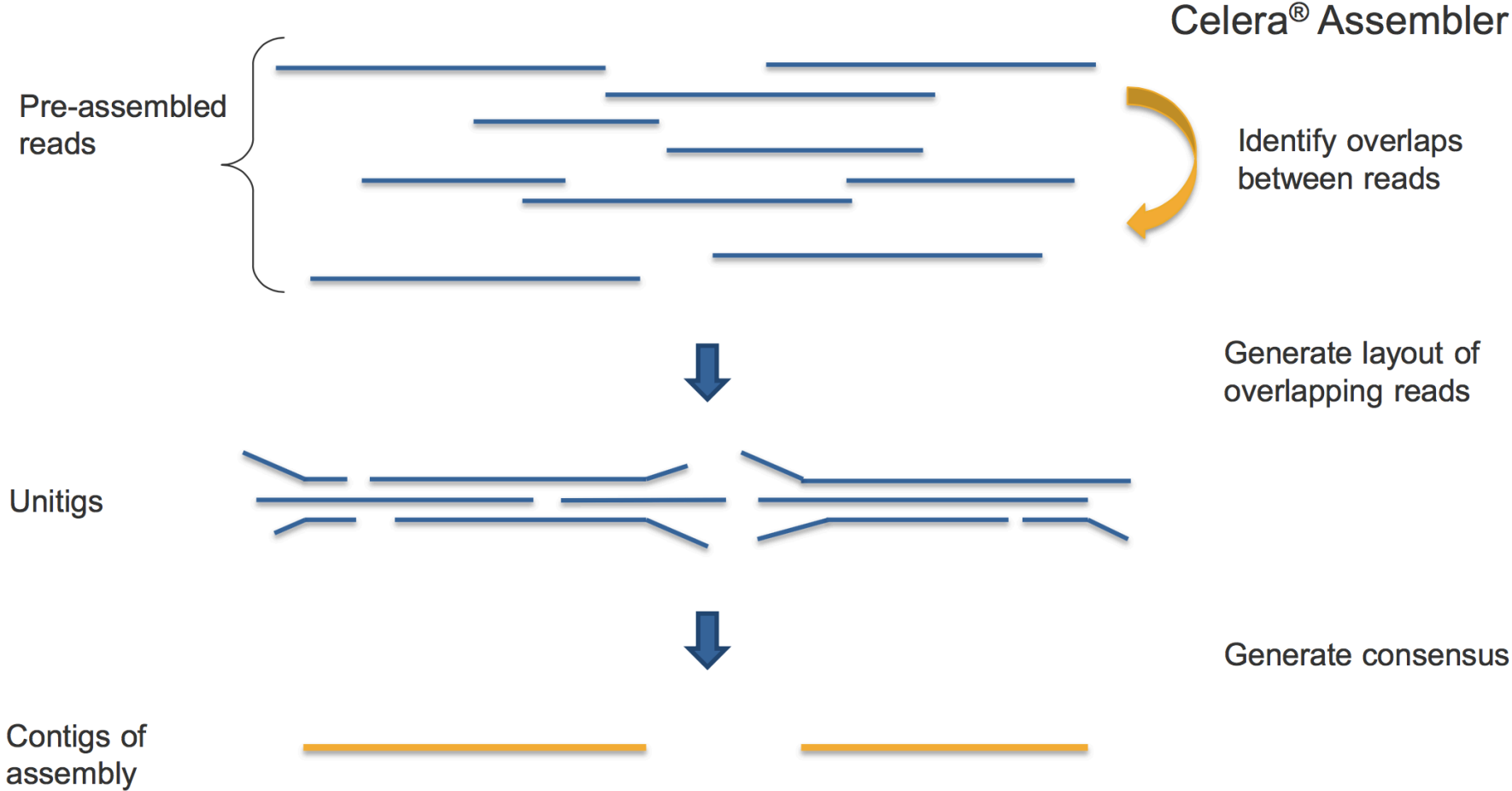
In this example, "T" insertion error corrected, and "CA" missing error recovered.

Supplementary Figure 12

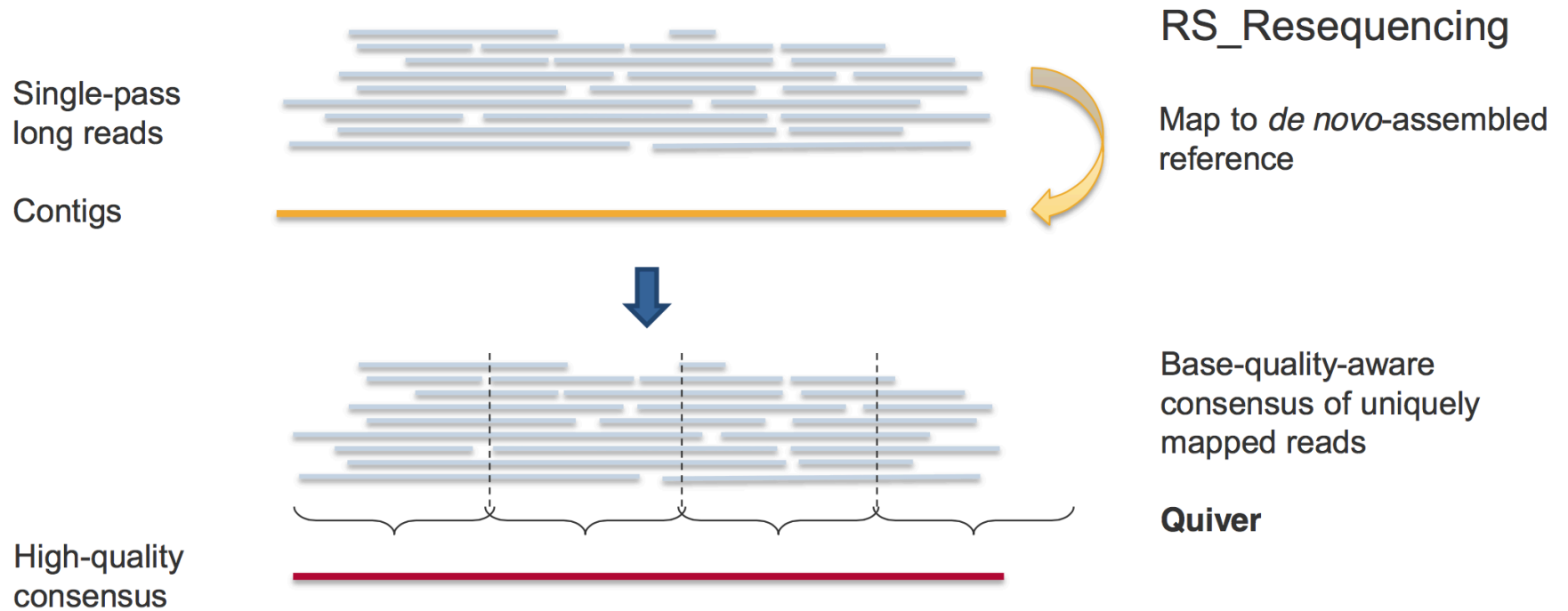
An Example of how the FALCON-sense algorithm generates consensus sequence.

- Popular long read assemblers:
  - HGAP (PacBio)
    - Limited to genomes < 200MB
    - <http://www.pacb.com/support/software-downloads/>
  - Canu (PacBio, Nanopore)
    - Any size genome
    - <https://github.com/marbl/canu>
  - Falcon (PacBio, Nanopore?)
    - Small and medium sized genomes, Diploid aware
    - <http://pb-falcon.readthedocs.io/en/latest/>
  - Miniasm (PacBio, Nanopore)
    - Any size genomes
    - <https://github.com/lh3/miniasm>









SMRT® Portal Home Admin Help About Welcome, administrator! Account Log Off

**DESIGN JOB**      **MONITOR JOBS**      **VIEW DATA**

 **Open Existing**

 **Create New**

 **Import and Manage**

**RECENT JOBS**

Job Name	Protocol	Reference Sequence	Started	Status	User
ugmExampleA_reseq	RS_Resequencing.1	ugmExampleA		Completed	kluong
NHGRI10_4cells_tandemAr	HGAP_Assembly_Artifact.1			Completed	kluong

## Import and Manage



### Manage Protocols

Create and edit standard protocols for secondary analysis jobs in SMRT Portal.

### Manage Reference Sequences

Import and manage reference sequences for resequencing and visualization with SMRT View.

### Import SMRT Cells

Import raw data from SMRT cells for analysis in SMRT Portal.

### Import SMRT Pipe Jobs

Import SMRT Pipe jobs for display in SMRT Portal.

### Import SMRT Cells

Import raw SMRT cell data for analysis in SMRT Portal.

#### File Paths

Searched locations:

- /mnt/data3/vol53/2450432
- /mnt/data3/vol53/2450434
- /mnt/data3/vol53/2450435
- /mnt/data3/vol53/2450436
- /mnt/data3/vol53/2450437

#### Enter A Folder Containing SMRT Cell Data

Path

- /mnt/data3/vol53/2450462
- /mnt/data3/vol53/2450464
- /mnt/data3/vol53/2450466
- /mnt/data3/vol53/2450469
- /mnt/data3/vol53/2450471
- /mnt/data3/vol53/2450472

## File Paths

### Searched locations:

/mnt/data3/vol53/2450432  
/mnt/data3/vol53/2450434  
/mnt/data3/vol53/2450435  
/mnt/data3/vol53/2450436  
/mnt/data3/vol53/2450437  
/mnt/data3/vol53/2450438

### Confirmation



Are you sure you want to scan the selected path(s)? It could take a long time.

OK

Cancel

/mnt/data3/vol53/2450464  
/mnt/data3/vol53/2450466  
/mnt/data3/vol53/2450469  
/mnt/data3/vol53/2450471  
/mnt/data3/vol53/2450472

SMRT® Portal Home Admin Help About Welcome, administrator! Account Log Off

**DESIGN JOB**      **MONITOR JOBS**      **VIEW DATA**

 **Open Existing**

 **Create New**

 **Import and Manage**

**RECENT JOBS**

Job Name	Protocol	Reference Sequence	Started	Status	User
ugmExampleA_reseq	RS_Resequencing.1	ugmExampleA		Completed	kluong
NHGRI10_4cells_tandemAr	HGAP_Assembly_Artifact.1			Completed	kluong

SMRT® Portal Home Tech Support Files Help About Welcome, nsisneros! Account Log Off

**DESIGN JOB** **MONITOR JOBS** **VIEW DATA**

Job Name  Comments  Groups all User

Protocols  Reference [None selected]

**SMRT Cells Available** (Viewing 1 - 28 of 28)

Sample	Uri
Seabury-11-20-091114	/mnt/data3/vol53/fas/da
Seabury-11-20-091114	/mnt/data3/vol53/fas/da
Seabury-11-20-091114	/mnt/data3/vol53/fas/da
Seabury-11-20-091114	/mnt/data3/vol53/fas/da
Laxiflora-090314	/mnt/data3/vol60/2420
Laxiflora-090314	/mnt/data3/vol60/2420
Laxiflora-081914	/mnt/data3/vol60/2420
Laxiflora-090314	/mnt/data3/vol60/2420
Laxiflora-081914	/mnt/data3/vol60/2420
Laxiflora-081914	/mnt/data3/vol60/2420
K.Laxiflora-PP-080714-.08nM-1xMB	/mnt/data3/vol60/2420
Laxiflora-090314	/mnt/data3/vol60/2420
K.Laxiflora-PP-080714-.08nM-1xMB	/mnt/data3/vol60/2420
K.Laxiflora-PP-080714-.08nM-5xMB	/mnt/data3/vol60/2420
K.Laxiflora-PP-080714-.08nM-1xMB	/mnt/data3/vol60/2420
Laxiflora-081914	/mnt/data3/vol60/2420
Laxiflora-081914	/mnt/data3/vol60/2420
K.Laxiflora-PP-080714-.08nM-5xMB	/mnt/data3/vol60/2420
K.Laxiflora-PP-080714-.08nM-1xMB	/mnt/data3/vol60/2420
Laxiflora-081914	/mnt/data3/vol60/2420
Laxiflora-090314	/mnt/data3/vol60/2420298/0044
Laxiflora-090314	/mnt/data3/vol60/2420298/0046
K.Laxiflora-PP-080714-.08nM-5xMB	/mnt/data3/vol60/2420298/0055

records to view)

Ve	User	Groups	Started	Uri
----	------	--------	---------	-----

Show/Hide Columns Group Rows Export Table Data

Start Save Copy Cancel

**Analysis**

Which type(s) of analysis would you like to perform on your sequencing data?

- Reference-based**  
Analyze PacBio sequence data and compare to a known reference sequence. (This Includes tasks such as resequencing, cDNA Q/C and mapping, minor variant detection, and base modification analysis.)
- De novo assembly**  
Assemble a genome from PacBio data, perform phasing on long amplicons such as the HLA region.
- Data Prep**  
Prepare PacBio sequencing data for analysis. This includes filtering data.
- Display all types of analysis**  
Show all available analysis types, including custom protocols. (A protocol is a set of software algorithms that performs analysis on your input sequencing data.)

Don't show this again

Next

SMRT® Portal Home Tech Support Files Help About Welcome, nsisneros! Account Log Off

**DESIGN JOB** **MONITOR JOBS** **VIEW DATA**

Job Name  Comments  Groups  User

Protocols  Reference

SMRT Cells in Job (No records to view)

Sample	Vol	User	Groups	Started	Uri
Seabury-11					3/fas/data_perm/inputs_dropbox/D06_3
Seabury-11					3/fas/data_perm/inputs_dropbox/D06_2
Seabury-11					3/fas/data_perm/inputs_dropbox/D06_4
Seabury-11					3/fas/data_perm/inputs_dropbox/D06_1
Laxiflora-09				0/2420298/0045	
Laxiflora-09				0/2420298/0048	
Laxiflora-08				0/2420298/0063	
Laxiflora-09				0/2420298/0042	
Laxiflora-08				0/2420298/0057	
Laxiflora-08				0/2420298/0059	
K.Laxiflora-09				0/2420298/0052	
Laxiflora-090314				/mnt/data3/vol60/2420298/0043	
K.Laxiflora-PP-080714-.08nM-1xMB				/mnt/data3/vol60/2420298/0051	
K.Laxiflora-PP-080714-.08nM-5xMB				/mnt/data3/vol60/2420298/0054	
K.Laxiflora-PP-080714-.08nM-1xMB				/mnt/data3/vol60/2420298/0049	
Laxiflora-081914				/mnt/data3/vol60/2420298/0064	
Laxiflora-081914				/mnt/data3/vol60/2420298/0061	
K.Laxiflora-PP-080714-.08nM-5xMB				/mnt/data3/vol60/2420298/0056	
K.Laxiflora-PP-080714-.08nM-1xMB				/mnt/data3/vol60/2420298/0050	
Laxiflora-081914				/mnt/data3/vol60/2420298/0062	
Laxiflora-090314				/mnt/data3/vol60/2420298/0044	
Laxiflora-090314				/mnt/data3/vol60/2420298/0046	
K.Laxiflora-PP-080714-.08nM-5xMB				/mnt/data3/vol60/2420298/0055	

Search

Show/Hide Columns Group Rows Export Table Data



SMRT® Portal Home Tech Support Files Help About Welcome, nsisneros! Account Log Off

**DESIGN JOB** **MONITOR JOBS** **VIEW DATA**

Job Name BFX course Comments Groups all User

Protocols RS\_HGAP\_Assembly.3

**SMRT Cells Available** (Viewing 1 - 28 of 28)

Sample	Uri
Seabury-11-20-091114	/mnt/d
Seabury-11-20-091114	/mnt/d
Seabury-11-20-091114	/mnt/d
Seabury-11-20-091114	/mnt/d
Laxiflora-090314	/mnt/d
Laxiflora-090314	/mnt/d
Laxiflora-081914	/mnt/d
Laxiflora-090314	/mnt/d
Laxiflora-081914	/mnt/d
Laxiflora-081914	/mnt/d
K.Laxiflora-PP-080714-.08nM-1xMB	/mnt/d
Laxiflora-090314	/mnt/d
K.Laxiflora-PP-080714-.08nM-1xMB	/mnt/d
K.Laxiflora-PP-080714-.08nM-5xMB	/mnt/d
K.Laxiflora-PP-080714-.08nM-1xMB	/mnt/d
Laxiflora-081914	/mnt/d
Laxiflora-081914	/mnt/d
Laxiflora-081914	/mnt/d
K.Laxiflora-PP-080714-.08nM-5xMB	/mnt/d
K.Laxiflora-PP-080714-.08nM-1xMB	/mnt/d
Laxiflora-081914	/mnt/d
Laxiflora-090314	/mnt/d
Laxiflora-090314	/mnt/d
K.Laxiflora-PP-080714-.08nM-5xMB	/mnt/d

Show/Hide Columns Group Rows Export Table Data

**Protocol Details For Job BFX Course**

Protocol: PreAssembler v2

Filtering:  Compute Minimum Seed Read Length

Control Filtering: Minimum Seed Read Length\* 6000

Assembly: Number Of Seed Read Chunks 6

Mapping: Alignment Candidates Per Chunk 10

Consensus: Total Alignment Candidates 24

Minimum Coverage For Correction 6

BLASR Options (Advanced) -noSplitSubreads -minReadLength

**AssembleUnitig v1**

Genome Size (Bp)\* 5000000

Target Coverage 25

Overlapper Error Rate 0.06

Overlapper Min Length 40

Overlapper K-mer 14

Pre-defined Spec File

OK Apply Cancel

**Minimum Seed Read Length:**

- 30X Coverage of longest Seed Reads automatically calculated
- Uncheck to override "auto"

**Key Parameter to set: Genome Size**

- 130 MB limit in SMRT Portal 2.3

Start Save Copy Cancel

SMRT® Portal Home Tech Support Files Help About Welcome, smrtuser! Account Log Off

**DESIGN JOB** **MONITOR JOBS** **VIEW DATA**

Job Name  Comments  Groups  User

Protocol  Reference  ✔ Completed

SMRT Cells: 1 Movies: 1

**REPORTS**

**GENERAL**

- Overview
- Filtering
- Subread Filtering

**DIAGNOSTIC**

- Adapters
- Loading

**RESEQUENCING**

- Mapping
- Coverage

**ASSEMBLY**

- Pre-Assembly
- Polished Assembly

**DATA**

**GENERAL**

- SMRT Cells in Job
- Reads [H5 FASTA](#)
- Filtering [CSV](#)
- Filtered Subreads [CSV FASTA FASTQ](#)

**ASSEMBLY**

- Preassembled Reads [FASTA FASTQ](#)
- Polished Assembly [FASTQ CSV FASTA](#)

**RESEQUENCING**

- Aligned Reads [H5 BAI BAM SAM](#)

Job Metric	Value
Polished Contigs	1
Adapter Dimers (<100bp)	0.00%
Short Inserts (11-100bp)	0.0%
Number of Bases	503,112,125
Number of Reads	59,211
N50 Read Length	12,848
Mean Read Length	8,496
Mean Read Score	0.84
Mapped Reads	56,207
Mapped Read Length of Insert	7,363
Average Reference Length	4,669,315
Average Reference Bases Called	100.0%
Average Reference Consensus Concordance	99.98%
Average Reference Coverage	94.12

**Adapters**

Observed Insert Length Distribution Histogram

**Subread Filtering**

Subread Filtering

**Mapping**

Mapped Subread Concordance

**Mapping**

Mapped Subread Length

**Mapping**

Mapped Polymerase Read Length

**Coverage**

Coverage Across Reference

**Coverage**

**Corrections**

- Running HGAP (Command line)
  - Install SMRT Analysis software
  - Make a HGAP assembly job using the SMRT portal and save.
  - Save the settings.xml file as HGAP\_protocol.xml
    - Every SMRT Portal job has the following structure. **Example:**

```
/path/to/smrtanalysis/userdata/jobs/016/016234
├─ data/
├─ results/
├─ log/
├─ workflow/
├─ job.sh
├─ input.xml
└─ settings.xml
```

- `data` is a **directory** that contains intermediate and final data files for the analysis job
- `results` is a **directory** that contains summary statistics and plots for the analysis job
- `log` is a **directory** that contains all log files for the analysis job
- `workflow` is a **directory** that contains all the executables for the analysis job
- `job.sh` is an executable file used by SMRT Portal to run the `smrtpipe.py` analysis job
- `input.xml` is a .xml file containing a list of input `bax.h5` files used to run the analysis job
- `settings.xml` is a .xml file containing the parameters needed to perform the analysis job

- Running HGAP (Command line) cont'd.
  - Modify Genome size in HGAP\_protocol.xml
    - `<param name="genomeSize" label="Genome Size (bp)">  
<value>5000000</value>`
  - Source the SMRT analysis environment
    - `source /path/to/smrtanalysis/install/smrtanalysis_2.3.0.140936/  
etc/setup.sh`
  - Add the full paths of your raw data (\*.bax.h5) into an input.fofn
    - `find <data_dir> -name "*.bax.h5" > input.fofn`
  - Convert the input.fofn to an input.xml
    - `fofnToSmrtpipeInput.py input.fofn > input.xml`
  - Run SMRT pipe using the protocol and input xmls.
    - `smrtpipe.py --params=HGAP_protocol.xml xml:input.xml`
  - Results are found in index.html in the working directory
  - Assembly is in data/polished\_assembly.fastq.gz

```
#!/bin/bash
#SBATCH -A <your uppmax project>
#SBATCH -p core
#SBATCH -n 8
#SBATCH -t 1-00:00:00
#SBATCH -J run_smrt_assembly
#SBATCH -e run_smrt_assembly-%j.out
#SBATCH -o run_smrt_assembly-%j.out

module load bioinfo-tools SMRT/2.3.0
WORK_DIR=$SNIC_TMP/smrt_assembly_$(date +%Y_%m_%d-%H.%M)
PROJ_DIR=$PWD
PROTOCOL_XML=$PROJ_DIR/Settings/HGAP_protocol.xml
DATA_DIR=${PROJ_DIR}/00_RawData          # Use full path
GENOME_SIZE=5000000

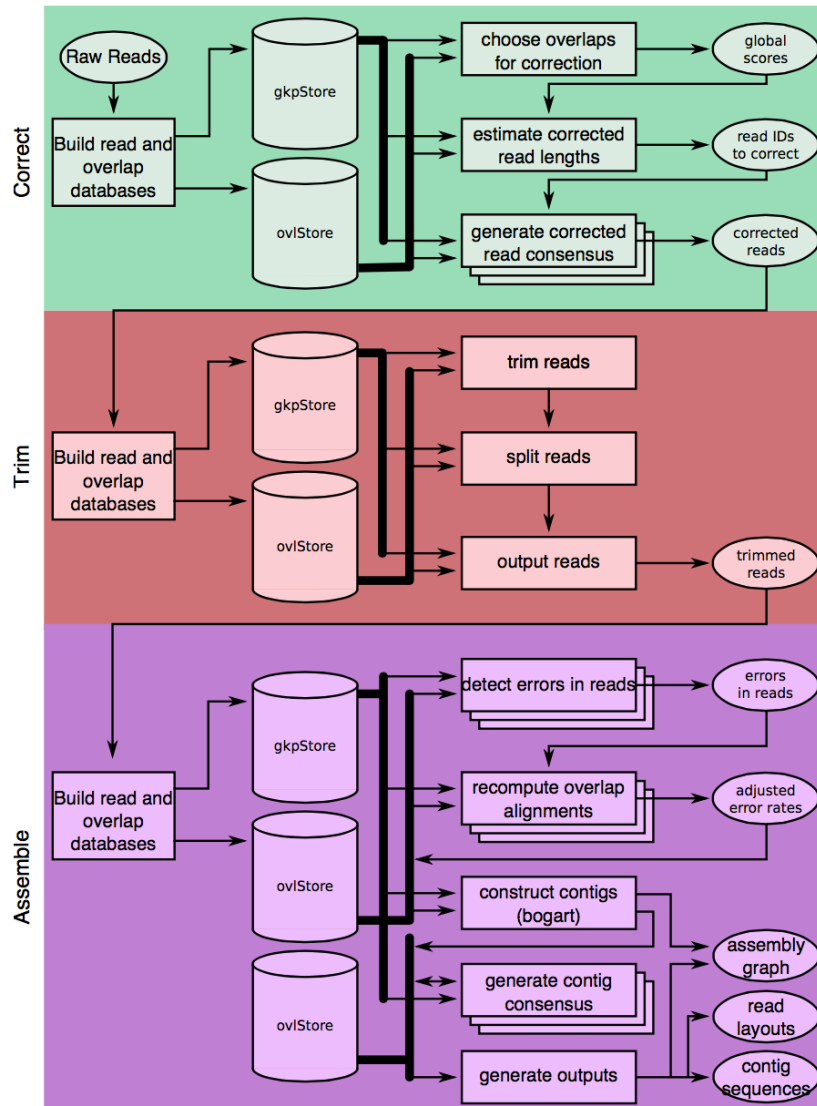
# Modify Protocol xml to the correct genome size
perl -0777 -i.original -pe "s/<param name=\"genomeSize\" label=\"Genome Size \\\(bp\\\\)\">\\n\\s+<value>\\d
+<\\/value>/<param name=\"genomeSize\" label=\"Genome Size (bp)\">\\n\\t\\t<value>$GENOME_SIZE<\\/value>/
igs" $PROTOCOL_XML

# Activate SMRT Analysis environment
source $SMRT_SETUP_SCRIPT
mkdir -p $WORK_DIR; cd $WORK_DIR

# Make input file
find ${DATA_DIR} -name "*.bax.h5" > input.fofn
fofnToSmrtpipeInput.py input.fofn > input.xml

smrtpipe.py --params=$PROTOCOL_XML xml:input.xml

cd $PROJ_DIR; rsync -av $WORK_DIR .
```



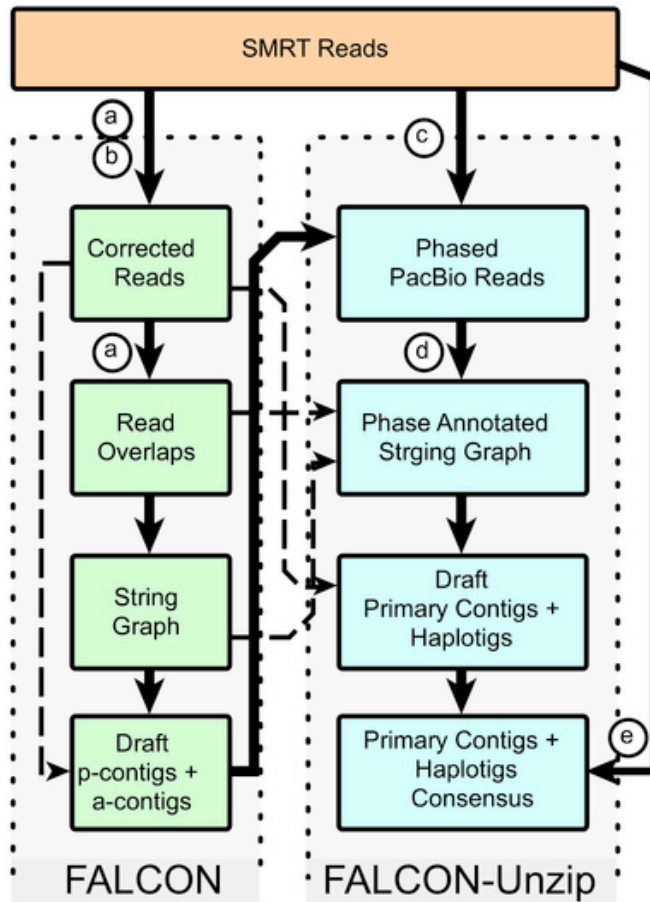
**Figure 1. A full Canu run includes three stages: correction (green), trimming (red), and assembly (purple).** Canu stages share an interface for binary on-disk stores (databases) as well as parallel store construction. In all stages, the first step constructs an indexed store of input sequences, generates a k-mer histogram, constructs an indexed store of all-vs-all overlaps, and collates summary statistics. The correction stage (green) selects the best overlaps to use for correction, estimates corrected read lengths, and generates corrected reads. The trimming stage (red) identifies unsupported regions in the input and trims or splits reads to their longest supported range. The assembly stage (purple) makes a final pass to identify sequencing errors; constructs the best overlap graph; and outputs contigs, an assembly graph, and summary statistics.

- Running Canu
  - Can auto-detect cluster settings (not recommended for milou)
  - Run canu on a node
    - useGrid=false
    - maxThreads=\$NPROCS

```
canu -p <file_prefix> -d <out_dir> genomeSize="18m"  
maxThreads=24 useGrid=false -pacbio-raw  
<filtered_subreads.fastq.gz>
```

- Results
  - Sequence is in file\_prefix.contigs.fasta
  - Assembly graph is in file\_prefix.gfa

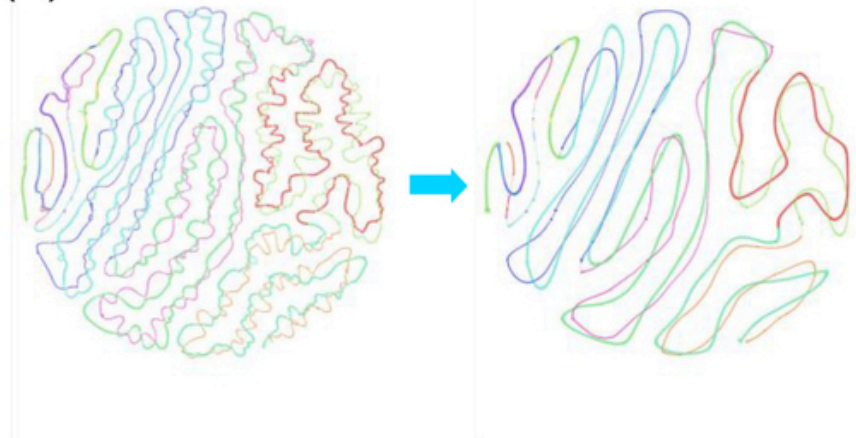
(a)



External code and internal modules used in FALCON and FALCON-Unzip

- (a) Daligner
- (b) Consensus Module (FALCON-sense)
- (c) Phasing Module (FALCON-phasing)
- (d) Graph "Unzip" Module
- (e) BLASR Alignment+ Quiver Consensus Module

(b)





## The fragment assembly string graph

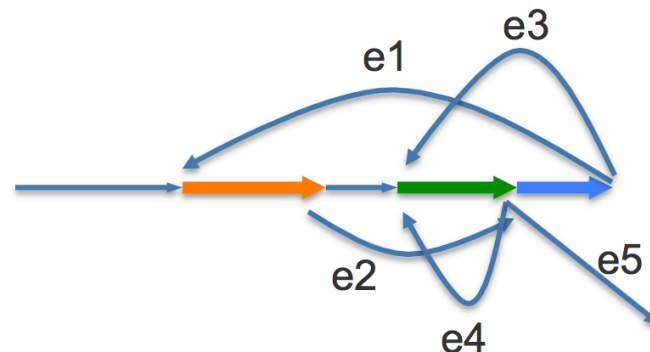
Eugene W. Myers

Department of Computer Science, University of California, Berkeley, CA, USA



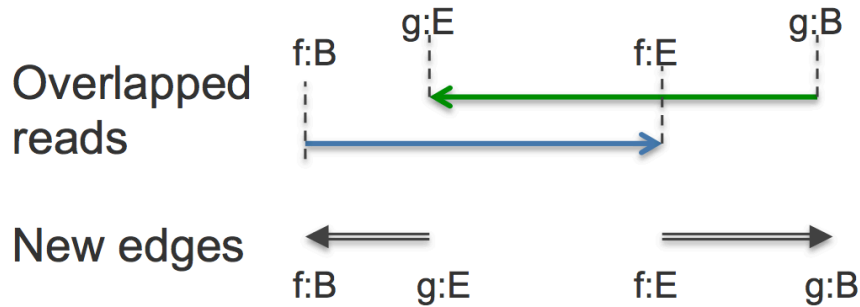
- String graph:
  - A graph structure that models a genome
- Nodes:
  - Particular positions (typically corresponding to the beginnings or endings of the read fragments) in the genome
- Edges:
  - The sequence between the vertices
- Any string from a path spell out a possible assembly from the reads

## String Graph

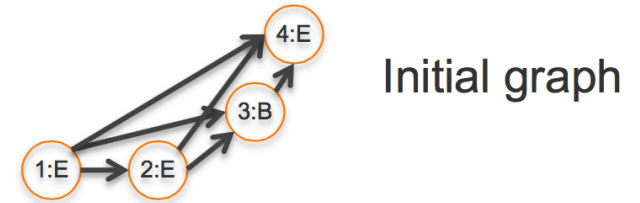
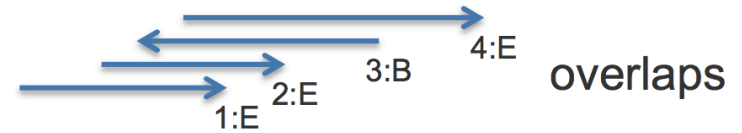


For each overlap, two edges are constructed.

Example:

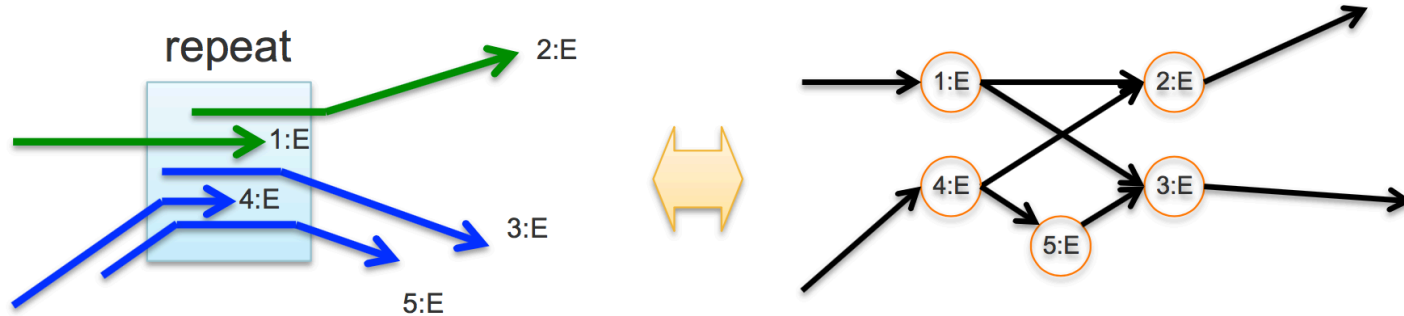


Add f:B, g:B, f:E, g:E as vertices  
 Add edges f:E → g:B and g:E → f:B



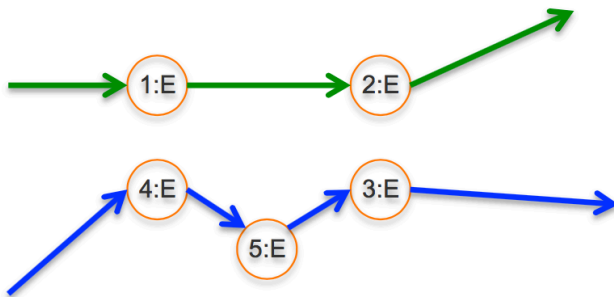
Transitive Reduction



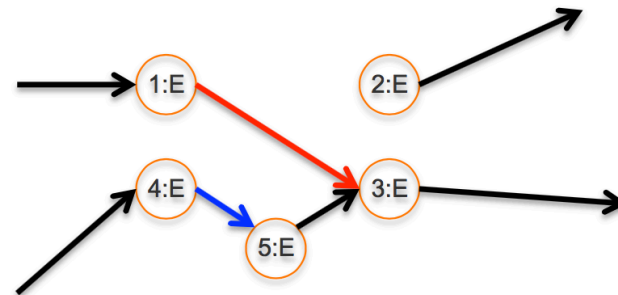


- Using a simple “best overlapping logic” to “untangle” the knots.

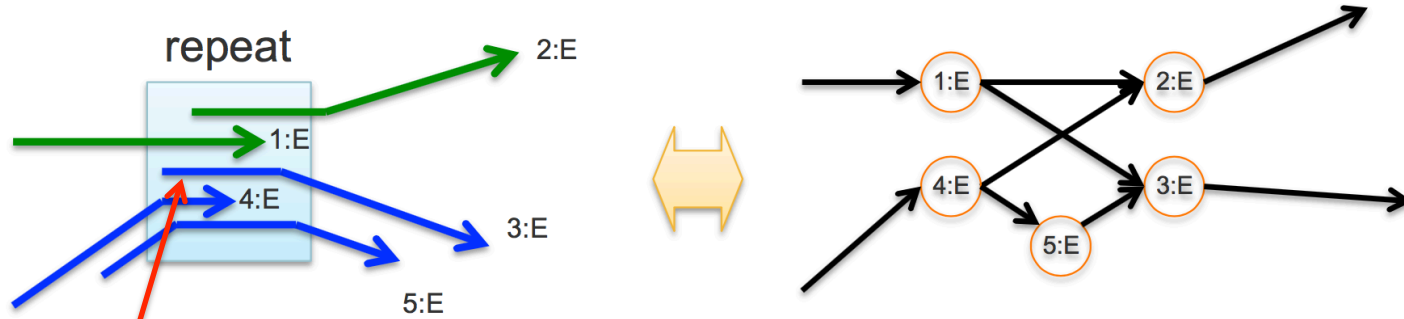
Desired final graph



Best overlap string graph

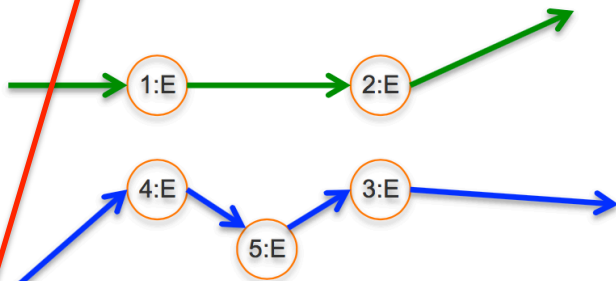


The 4E→5E edge is better than 4E →2E.  
 The 1E→3E edge is better than 1E →2E.  
 (“wrong” edge)

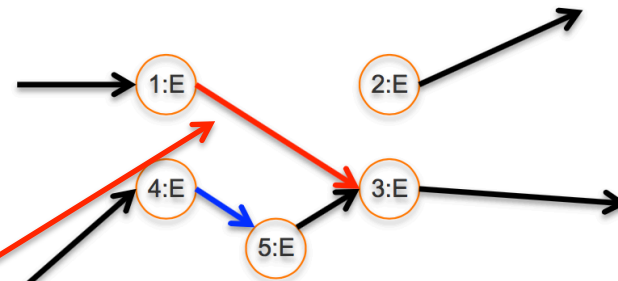


- Using a simple “best overlapping logic” to “untangle” the knots.

Desired final graph



Best overlap string graph



The 4E → 5E edge is better than 4E → 2E.  
 The 1E → 3E edge is better than 1E → 2E.  
 (“wrong” edge)

Wrong edge selected because of greater overlap with other read

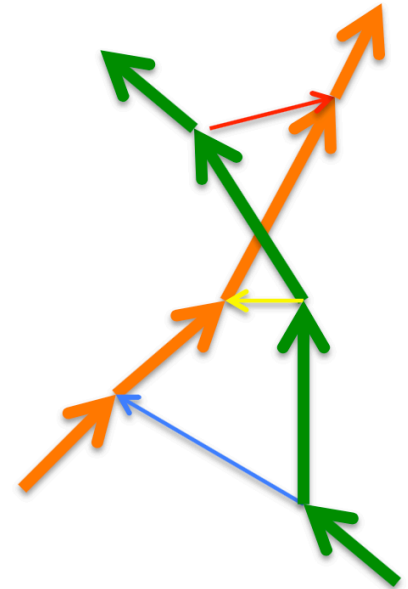
String Graph



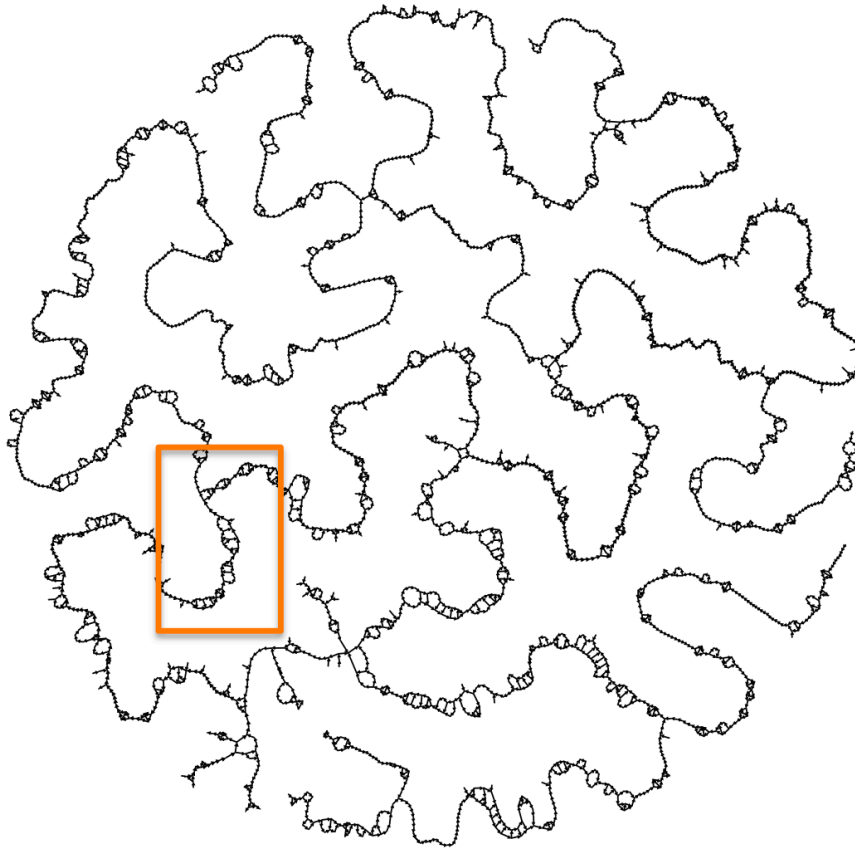
Unitig Graph



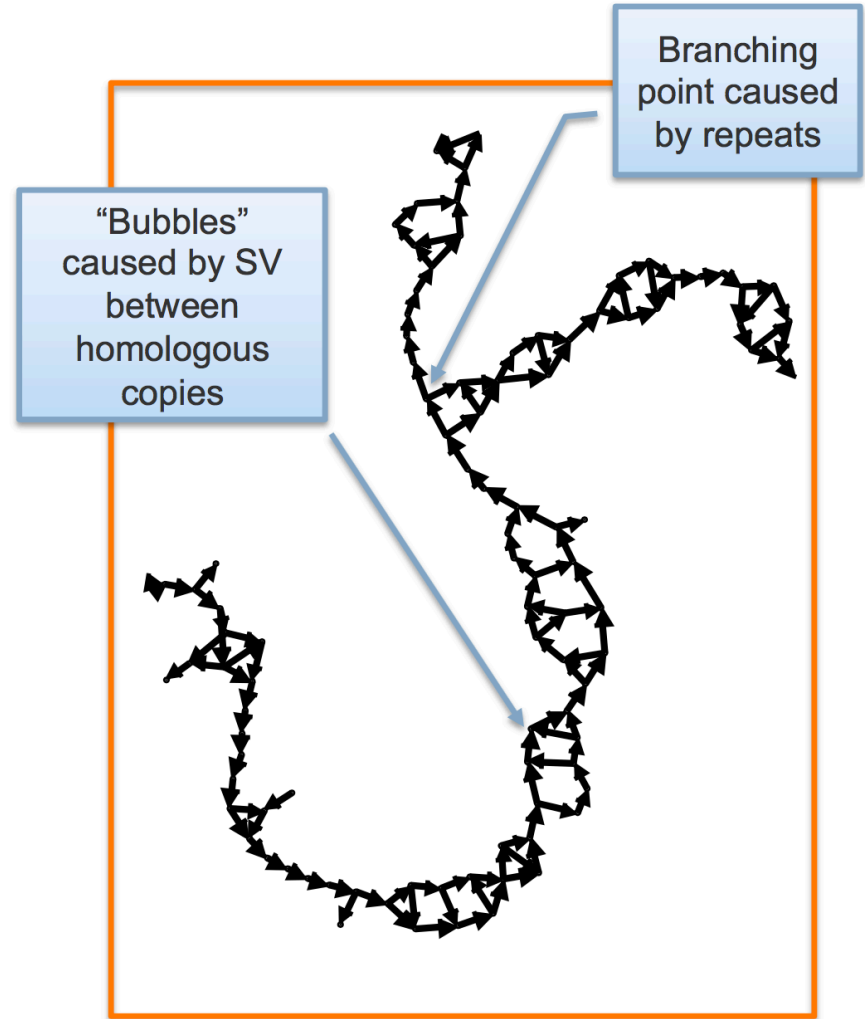
Graph traversal for generating contigs

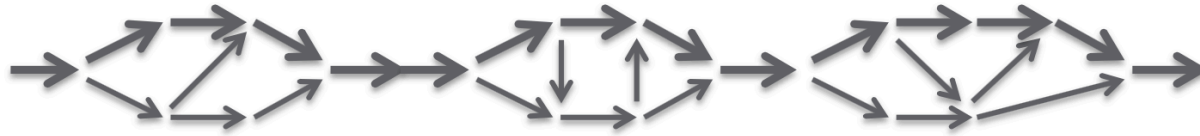


An unitig graph from Ler-0 + Col-0 data

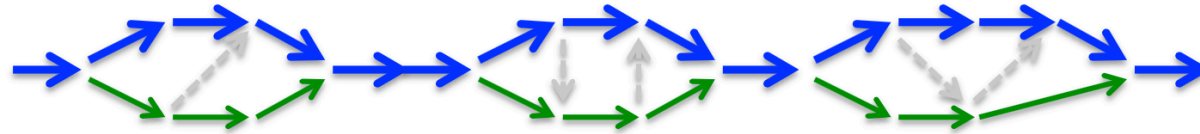


The graph "diameter" ~ 12 M bp  
Mean edge size=17.4 k bp

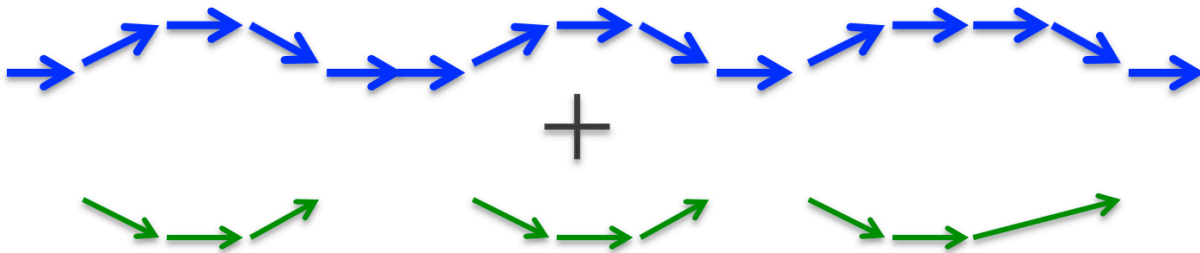




String Bundle



Choose a path to be the "primary contig"



Identify "associated contigs"

## Biological sequence

Two haplotypes with various degree of heterozygosity



The high heterozygosity part (orange) has no overlap with the counter part and becomes its own p-contig



## FALCON assembly model

The low heterozygosity part (purple) has no associated haplotig as no reliable phasing information available



## FALCON-Unzip assembly model



- Running Falcon
  - Make a configuration file
  - Can use SGE queuing manager
  - Run locally on a node
  - Separate filtered subreads into separate fasta's for each movie
    - `zcat *.fastq.gz | seqtk seq -l 5000 -A - | awk 'BEGIN { RS=">"; FS="/" } { print ">"substr($0,1,length($0)-1) > $1".fasta" }'`
  - Make an input fofn
    - `/bin/ls -l *.fasta > input.fofn`
  - Run Falcon
    - `fc_run.py falcon.cfg`

- Notes from the author
  - Falcon is limited by file i/o capabilities
    - Lustre file system recommended
    - NFS can handle 3-5 concurrent jobs during pre-assembly
    - Highly repetitive genomes require quadratically more storage space
  - Falcon scales quadratically
    - All-by-all comparison of raw subreads, with matches written to disk

- The Falcon config file (parameter rich - explanation in the exercises)

```
[General]
```

```
jobtype = local           # other values sge, slurm
```

```
input_fofn = input.fofn
```

```
input_type = raw         # uncorrected reads
```

```
#input_type = preads     # falcon corrected reads
```

```
# The length cutoff used for seed reads used in initial  
mapping - these make the corrected reads
```

```
length_cutoff = 12000    # use longest 30X coverage
```

```
# The length cutoff used for seed reads used for pre-  
assembly - the min length of corrected reads
```

```
length_cutoff_pr = 12000 # 0-5000 lower than above
```

- No read correction step
- Miniasm implements Overlap - Layout, Racon implements Consensus

```
# Overlap
```

```
minimap/minimap -Sw5 -L100 -m0 -t8 reads.fq reads.fq | \  
gzip -1 > reads.paf.gz
```

```
# Layout
```

```
miniasm/miniasm -f reads.fq reads.paf.gz > contigs.gfa
```

```
# Get fasta
```

```
awk '/^S/{print ">"+seq"\n"$3}' contigs.gfa > contigs.fasta
```

```
# Consensus
```

```
miniasm/minimap -t 16 contigs.fasta reads.fq | \  
racon -t 16 reads.fq - contigs.fasta consensus_contigs.fasta
```

- ABruijn
  - Uncorrected overlap assembly of long read sequences followed by polishing
  - <https://github.com/fenderglass/ABruijn>
- Ra
  - Uncorrected overlap assembly of long read sequences
  - <https://github.com/mariokostelac/ra-integrate>
- ARacon
  - Combination of GraphMap + Miniasm + Racon
  - <https://github.com/isovic/aracon>
- Hinge
  - Read filtering (but no correction) followed by overlap assembly of long read sequences
  - <https://github.com/fxia22/HINGE>
- SMARTdenovo
  - Uncorrected overlap assembly of long read sequences
  - <https://github.com/ruanjue/smarddenovo>

- Assembly Size
  - Assemblathon Script (<https://github.com/KorfLab/Assemblathon>)
  - Quast

Number of scaffolds	556	
Total size of scaffolds	31318563	← Expected genome size
Longest scaffold	447934	
Shortest scaffold	8580	
Number of scaffolds > 1K nt	556	100.0%
Number of scaffolds > 10K nt	555	99.8%
Number of scaffolds > 100K nt	38	6.8%
Number of scaffolds > 1M nt	0	0.0%
Number of scaffolds > 10M nt	0	0.0%
Mean scaffold size	56328	
Median scaffold size	43995	
N50 scaffold length	60037	
L50 scaffold count	152	

- Corrected Read Coverage
  - What happened in the correction process
  - High coverage? Use the ~100X longest subreads

```
-- Found 87386 reads.
-- Found 1654383605 bases (45.95 times coverage).
--
-- Read length histogram (one '*' equals 265.11 reads):
--      0      999      0
--     1000    1999      0
--     2000    2999      0
--     3000    3999      0
--     4000    4999      0
--     5000    5999      0
--     6000    6999      0
--     7000    7999      0
--     8000    8999      0
--     9000    9999      0
--    10000  10999      0
--    11000  11999      0
--    12000  12999      0
--    13000  13999      0
--    14000  14999      0
--    15000  15999  18558 *****
--    16000  16999  15099 *****
--    17000  17999  11974 *****
--    18000  18999   9486 *****
--    19000  19999   7344 *****
--    20000  20999   5652 *****
--    21000  21999   4328 *****
--    22000  22999   3516 *****
--    23000  23999   2725 *****
--    24000  24999   2057 *****
--    25000  25999   1672 *****
--    26000  26999   1243 ****
--    27000  27999    920 ***
--    28000  28999    735 **
--    29000  29999    541 **
--    30000  30999    414 *
--    31000  31999    324 *
```

- Falcon: DBstats 1-preads\_ovl/preads.db
  - Focus on % Bases column (multiply by read coverage to find cutoff).

Statistics for all wells of length 500 bases or more

12,915 reads out of 13,124 ( 98.4%)  
 116,202,931 base pairs out of 116,263,784 ( 99.9%)

8,997 average read length  
 6,983 standard deviation

Base composition: 0.249(A) 0.239(C) 0.258(G) 0.255(T)

Distribution of Read Lengths (Bin size = 1,000)

Bin:	Count	% Reads	% Bases	Average
42,000:	1	0.0	0.0	42279
41,000:	2	0.0	0.1	41631
...				
(more bin values)				
...				
3,000:	1,065	75.5	95.0	11317
2,000:	1,328	85.8	97.9	10259
1,000:	1,444	97.0	99.7	9251
0:	387	100.0	100.0	8997

Find the row that corresponds to 30x coverage.

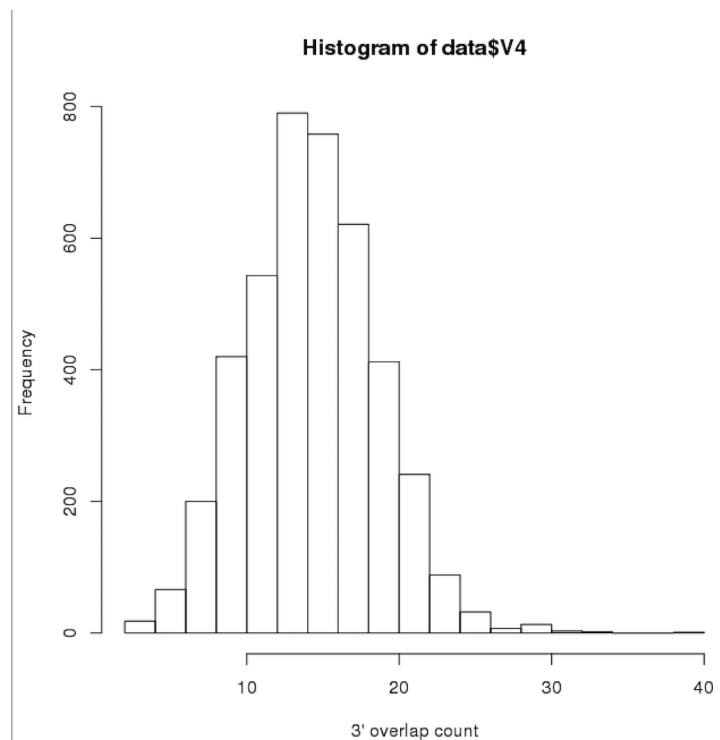
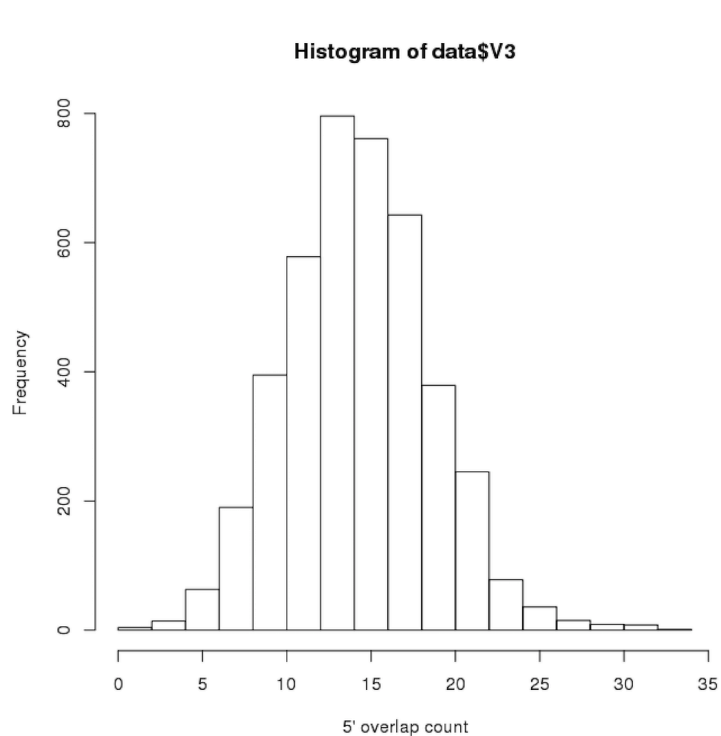
Total bases/ Est. genome size = Total coverage.  
 30x / Total coverage = % bases

Bin corresponding to 30x is the cutoff length.



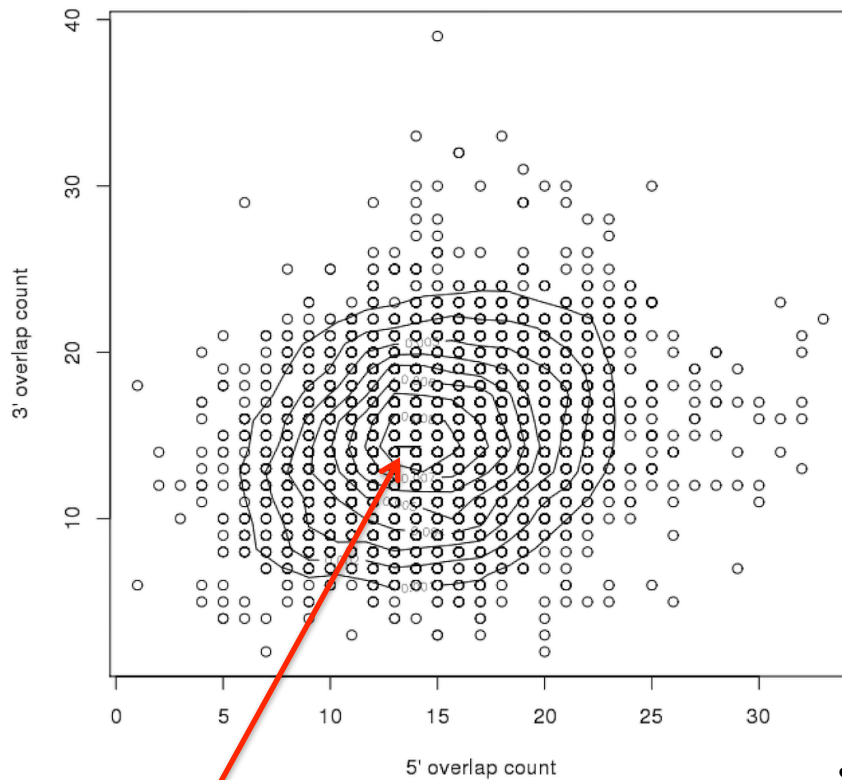
- Falcon: Overlap statistics

```
- cd 1-preads_ovl/ ;  
fc_ovlp_stats --fofn merge-gather/las.fofn >  
ovlp_stats.txt
```

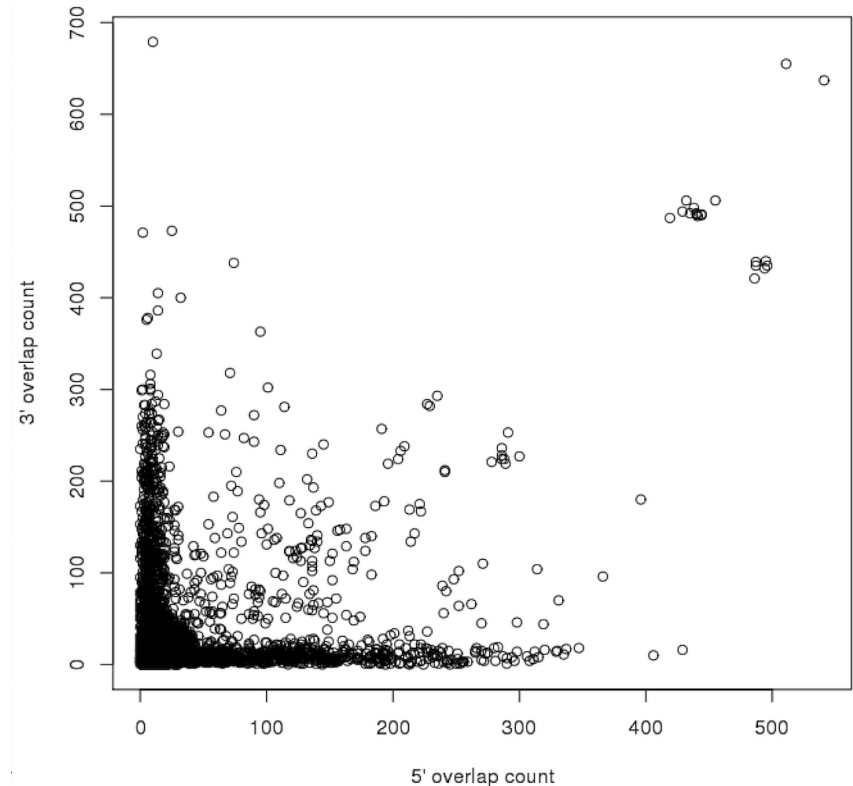


```
$ R  
> data <- read.table("ovlp_stats.txt")  
> hist(data$V3,xlab="5' overlap")  
> hist(data$V4,xlab="3' overlap")
```

- Falcon: Overlap statistics



Central density and symmetry along  $x=y$  is expected.



```
$  
> library(MASS)  
> data <- read.table("ovlp_stats.txt")  
> plot(data$V3,data$V4)  
> z <- kde2d(data$V3,data$V4)  
> contour(z,add=TRUE)
```

- Assembly Graph
  - Check connectedness of contigs
    - Is longer range information needed?
      - Higher quality sequence material
      - BioNano
      - Chicago / Dovetail

Nodes: 8,046  
Edges: 341  
Total length: 57,099,691

Graph drawing

Scope: Entire graph

Style:  Single  Double

Draw graph

Graph display

Zoom: 10.6%

Node width: 25.0

Random colours

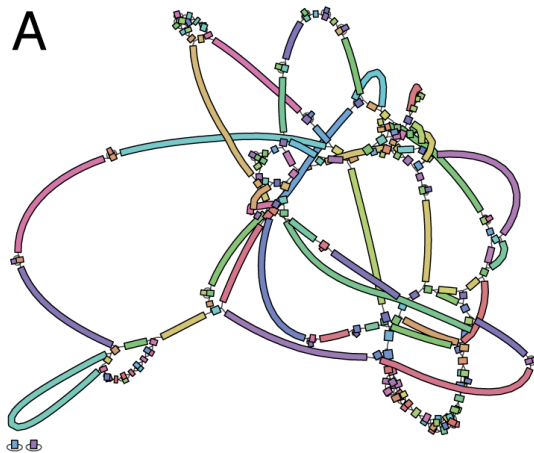
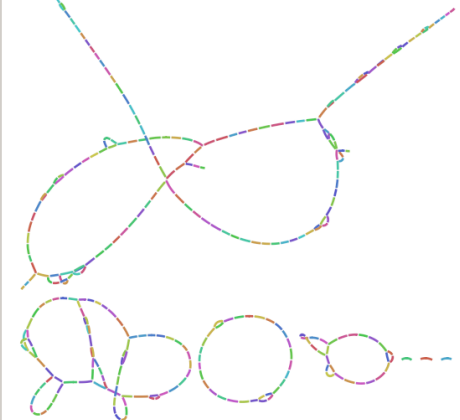
Node labels

Custom  Name

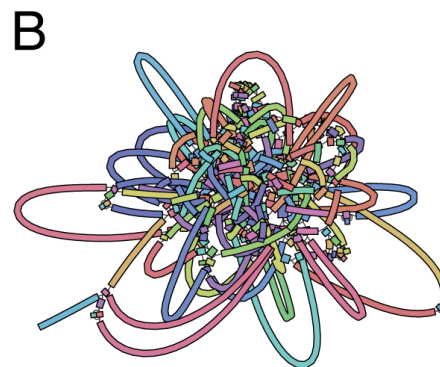
Length  Read depth

BLAST hits  CSV data

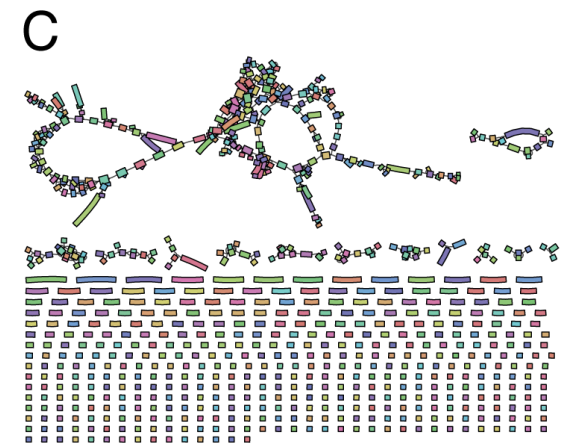
Font  Text outline



Good

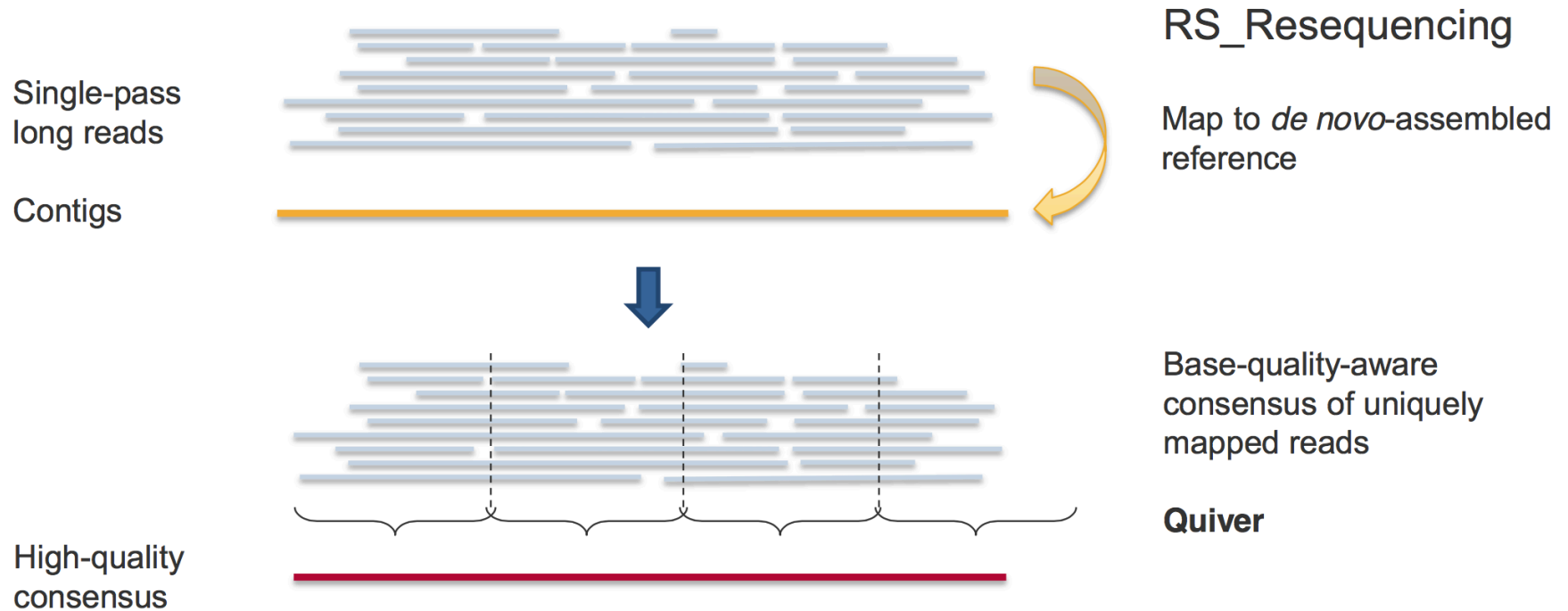


OK



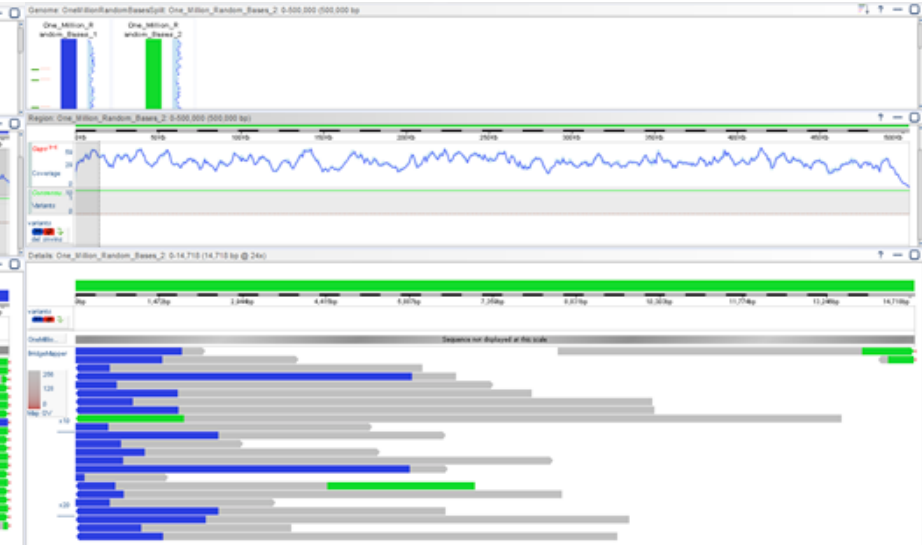
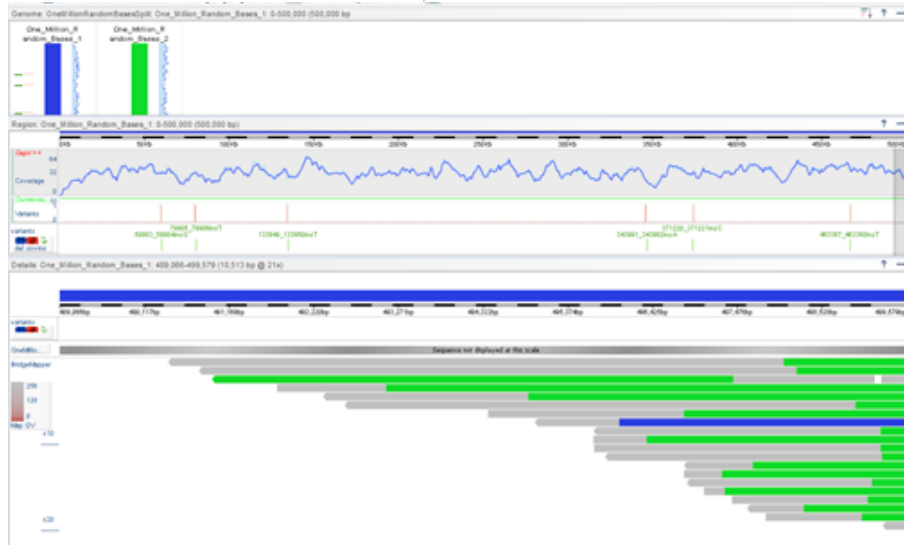
Bad

- Draft assemblies still contain many InDel and base substitution errors.
  - Correction using Quiver / Arrow and PacBio reads

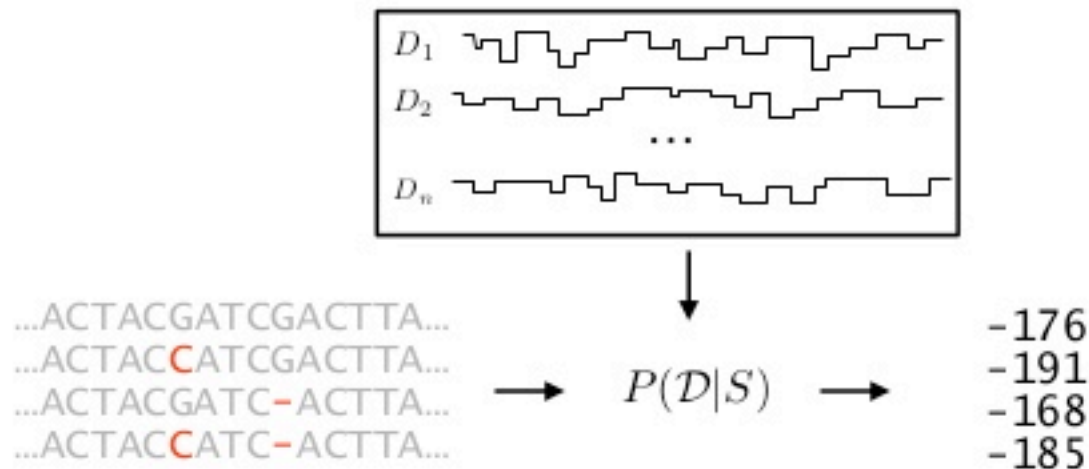


- In the SMRT portal select a protocol that includes the following modules
  - P\_Filter
  - P\_Mapping
  - P\_GenomicConsensus
    - RS\_Resequencing
    - RS\_BridgeMapper
      - P\_BridgeMapper
  
- To run via SMRT pipe
  - Copy the settings.xml of the dummy job
  - Make the draft assembly a reference using ReferenceUploader
  - Change reference value in settings.xml
  - Run

- Bridge Mapper results opened with SMRTview



- Draft assemblies still contain many InDel and base substitution errors.
  - Correction using Nanopolish and nanopore reads
    - Toolkit for working with signal-level data
    - Originally designed for improving a consensus sequence using the signals from multiple reads



## Human Assembly Consensus



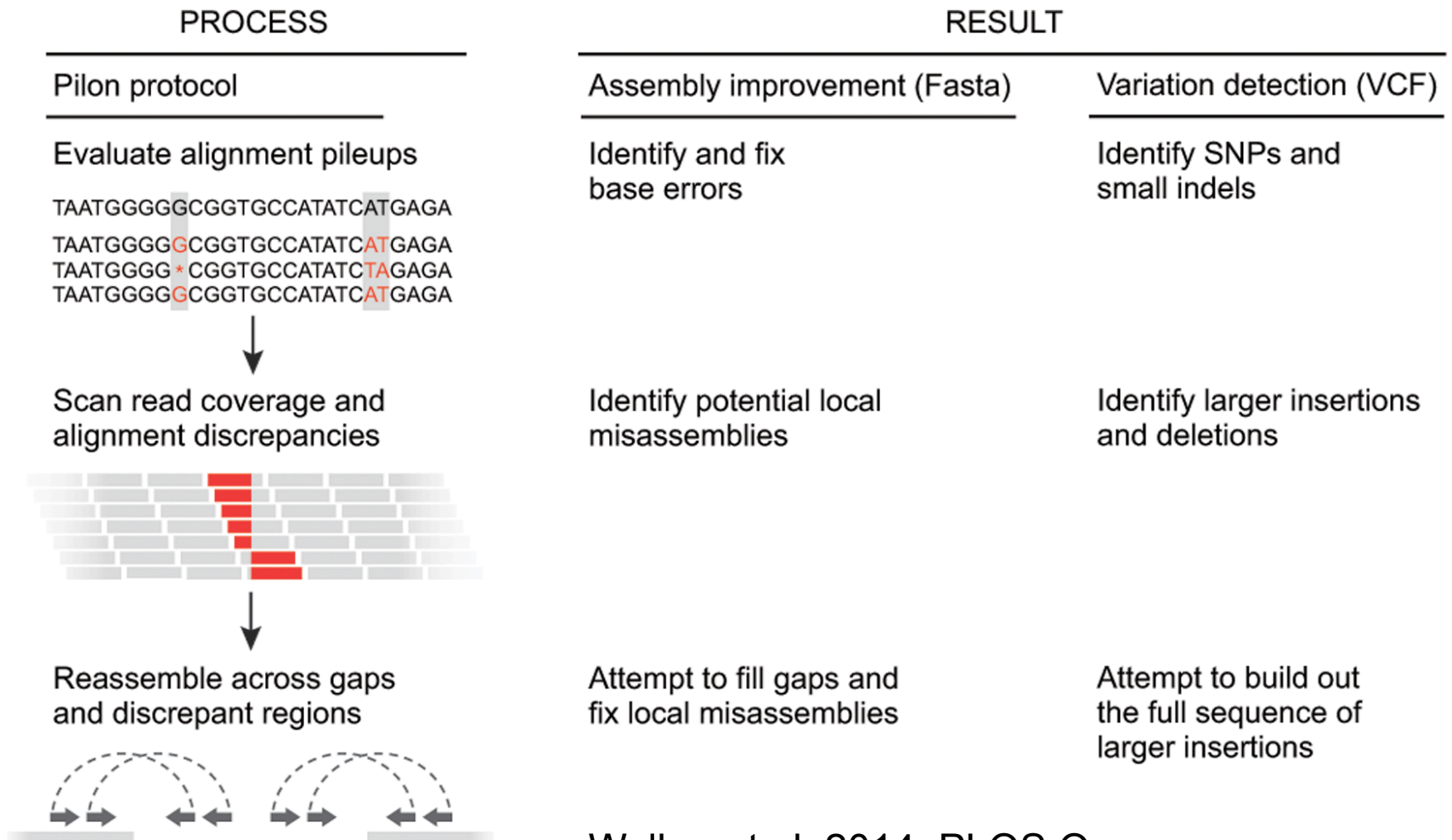
Assembly	Percent Identity
canu	94.8%
canu+racon	96.5%
canu+racon+nanopolish (30X)	99.1%
canu+racon+nanopolish (60X)	99.4%
canu+racon+nanopolish (60X) + pilon	99.6%

- This is a work in progress
- Polishing a single 6 Mbp chr20 contig
- 30X data set is NAI2878 consortium data only
- 60X data set includes 30X PCR-amplified NAI2878 provided by ONT
- stats calculated from bwa mem alignments to GRCh38
- differences that matched an NAI2878 variant were not consider an error

21



- Assembly Polishing
  - Can also be performed using Pilon and Illumina reads



- Long read sequencing is very dependent on sample DNA quality
- The longest reads are targeted for correction
- Correction with Illumina only does part of the read correction job
- Check basic stats
- Use the assembler output to diagnose major issues.
- Select your best assemblies
- Polish!
- Assess correctness.