

# ***De novo* meta-assembly of ultra-deep sequencing data**

Hamid Mirebrahim<sup>1</sup>, Timothy J. Close<sup>2</sup> and Stefano Lonardi<sup>1</sup>

<sup>1</sup>Department of Computer Science and Engineering

<sup>2</sup>Department of Botany and Plant Sciences



# Ultra-deep sequencing (>1,000x coverage) is possible and feasible, expected to become more common

HUMAN MUTATION Mutation in Brief 32: E1999-E2017 (2010) Online

## MUTATION IN BRIEF

HUMAN MUTATION

### Massive Parallel DNA Pyrosequencing Analysis of the Tumor Suppressor BRG1/SMARCA4 in Lung Primary Tumors



### Subclonal phylogenetic structures in cancer revealed by ultra-deep sequencing

Peter J. Campbell\*, Erin D. Pleasance\*, Philip J. Stephens\*, Ed Dicks\*, Richard Rance\*, Ian Goodhead\*, George A. Follows<sup>1</sup>, Anthony R. Green<sup>2</sup>, P. Andy Futreal<sup>1</sup>, and Michael R. Stratton\*<sup>1,5</sup>

\*Wellcome Trust Sanger Institute, Hinxton CB10 1SA, United Kingdom; <sup>1</sup>Department of Haematology, University of Cambridge, Cambridge CB2 2XY, United Kingdom; and <sup>2</sup>Institute of Cancer Research, Sutton, Surrey SW7 3RP, United Kingdom

Edited by Marshall Horwitz, University of Washington, Seattle, WA 98195, and accepted by the Editorial Board June 9, 2008 (received for review February 21, 2008)

During the clonal expansion of cancer from an ancestral cell with an initiating oncogenic mutation to symptomatic neoplasm, the occurrence of somatic mutations (both driver and passenger) can be used to track the on-going evolution of the neoplasm. All subclones within a cancer are phylogenetically related, with the prevalence of each subclone determined by its evolutionary fitness and the timing of its origin relative to other subclones. Recently developed, massively parallel sequencing platforms promise the

discovery of cells carrying drug resistance mutations before the initiation of therapy in both cancer (7) and infectious diseases (8, 9). To date, deep resequencing has detected variants down to a frequency of 1 in 100 (7–9), but its sensitivity for the detection of rarer variants has not been tested. With the appropriate informatic analyses and experimental design, the depth and breadth of sequencing available on the next-generation platforms will provide the tools to recon-

## Intervirology

### Original Paper

Intervirology 2014;57:384–392  
DOI: [10.1159/000368424](https://doi.org/10.1159/000368424)

Received February 19, 2014  
Accepted after revision August 5, 2014  
Published online October 31, 2014

### A Deep-Sequencing Method Detects Drug-Resistant Mutations in the Hepatitis B Virus in Indonesians

OPEN ACCESS Freely available online

### Transmission of Single HIV-1 Genomes and Dynamics of Early Immune Escape Revealed by Ultra-Deep Sequencing

Will Fischer<sup>1</sup>, Vitaly V. Ganusov<sup>1,2,3</sup>, Elena E. Giorgi<sup>1,3,3</sup>, Peter T. Hraber<sup>1,3</sup>, Brandon F. Keele<sup>4,5</sup>, Thomas Leitner<sup>1,3</sup>, Cliff S. Han<sup>1</sup>, Cheryl D. Gleasner<sup>1</sup>, Lance Green<sup>1</sup>, Chien-Chi Lo<sup>1</sup>, Ambarish Nag<sup>1</sup>, Timothy C. Wallstrom<sup>1</sup>, Shuyi Wang<sup>5</sup>, Andrew J. McMichael<sup>6</sup>, Barton F. Haynes<sup>7</sup>, Beatrice H. Hahn<sup>5</sup>, Alan S. Perelson<sup>1</sup>, Persephone Borrow<sup>8</sup>, George M. Shaw<sup>5</sup>, Tanmoy Bhattacharya<sup>1,9</sup>, Bette T. Korber<sup>1,9a</sup>

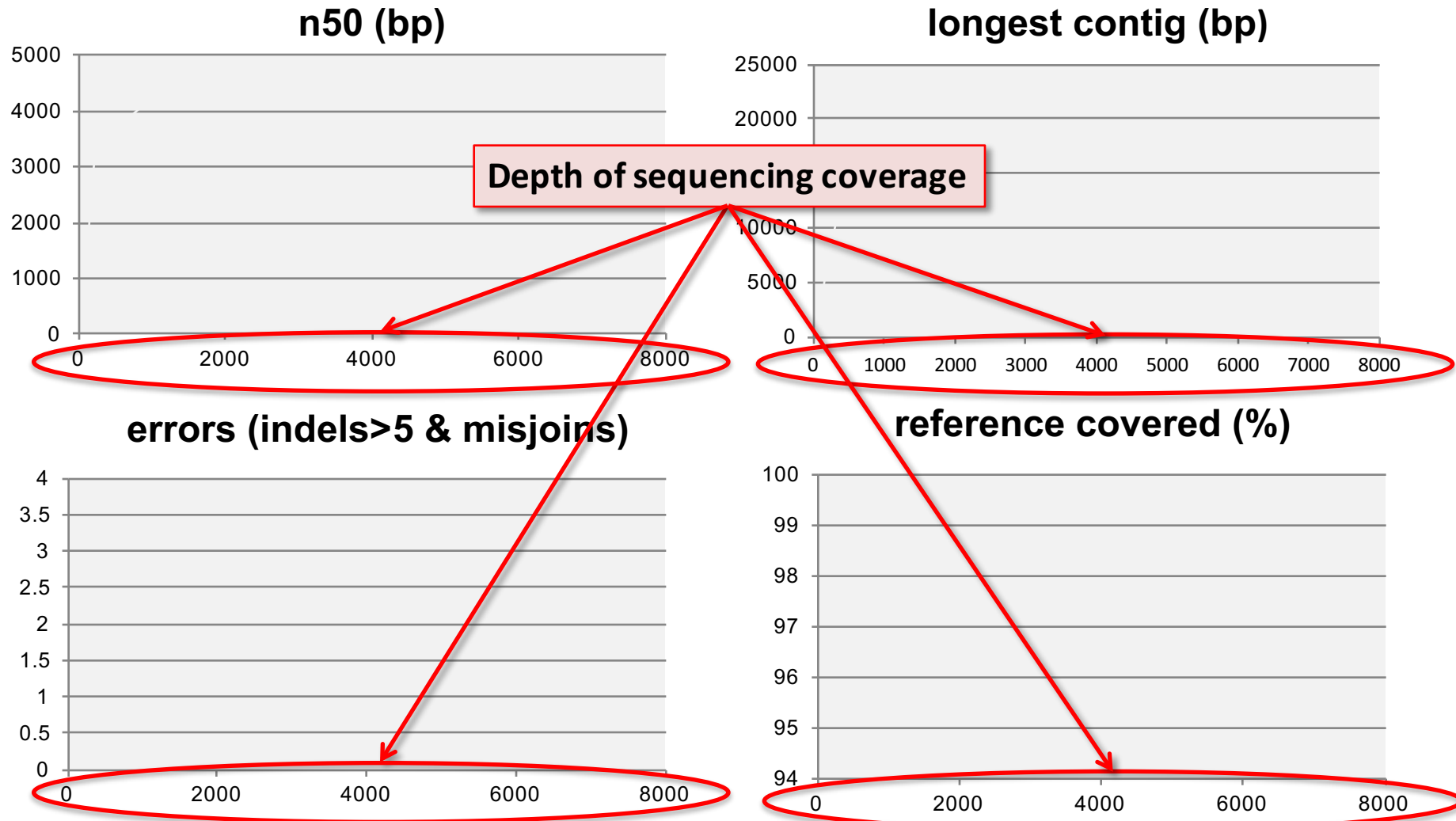
<sup>1</sup>Theoretical Biology, Los Alamos National Laboratory, Los Alamos, New Mexico, United States of America, <sup>2</sup>Department of Microbiology, University of Tennessee, Knoxville, Tennessee, United States of America, <sup>3</sup>Department of Mathematics and Statistics, University of Massachusetts, Amherst, Massachusetts, United States of America, <sup>4</sup>SAIC-Frederick, National Cancer Institute, Frederick, Maryland, United States of America, <sup>5</sup>Department of Medicine, University of Alabama at Birmingham, Birmingham, Alabama, United States of America, <sup>6</sup>Weatherall Institute of Molecular Medicine, Oxford University, Oxford, United Kingdom, <sup>7</sup>Duke University Medical Center, Durham, North Carolina, United States of America, <sup>8</sup>The Jenner Institute, University of Oxford, Compton, United Kingdom, <sup>9</sup>The Santa Fe Institute, Santa Fe, New Mexico, United States of America

### Intraclonal Diversity in Follicular Lymphoma Analyzed by Quantitative Ultradeep Sequencing of Noncoding Regions

...ence,<sup>‡</sup> and W. Richard Burack\*

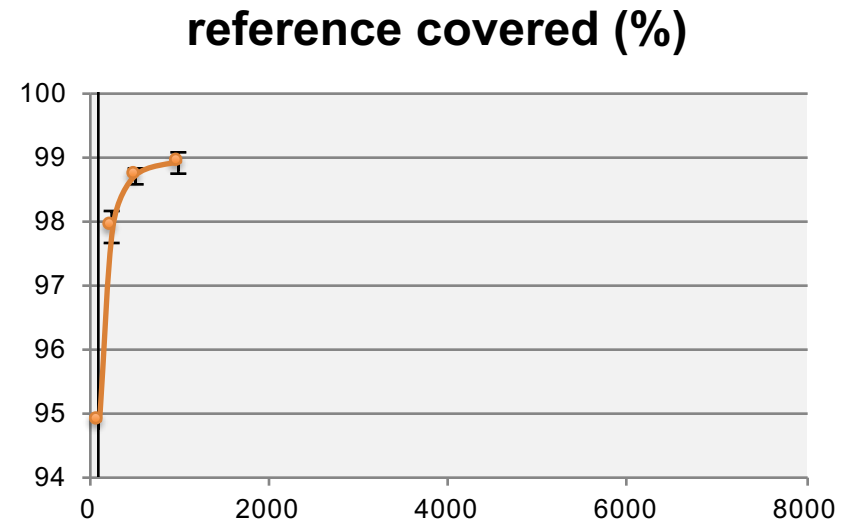
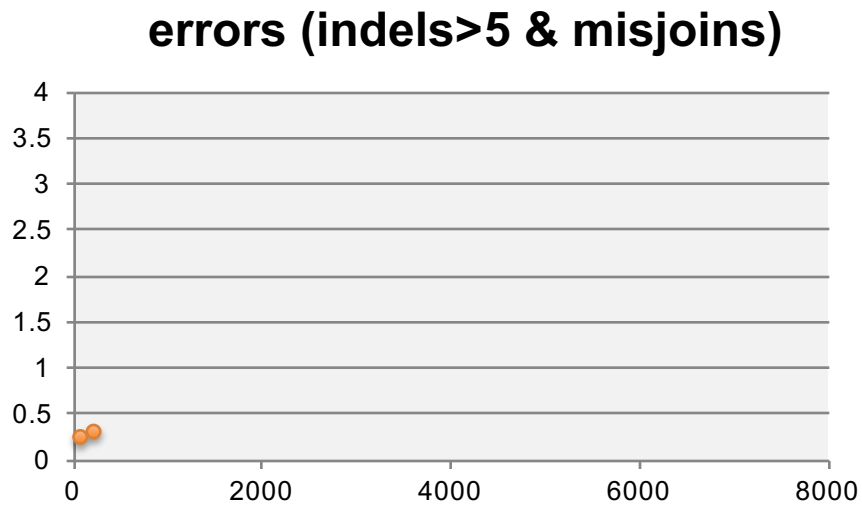
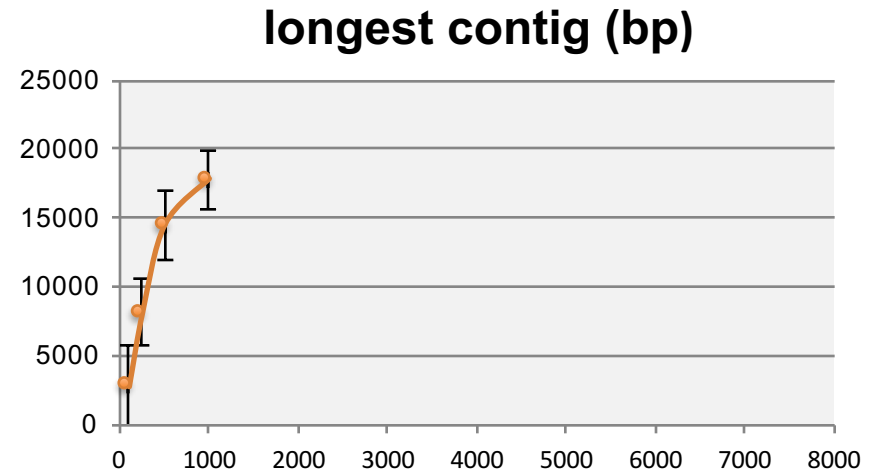
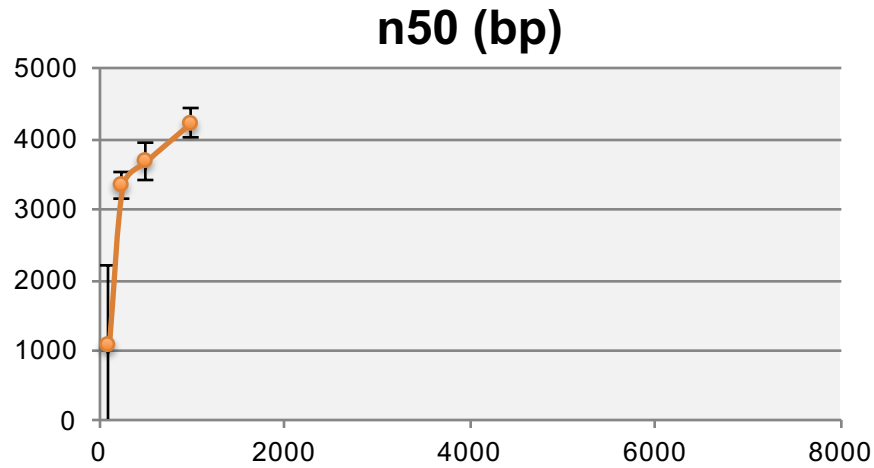
...intraclonal diversity is a prerequisite for tumor evolution. ... Intracloal heterogeneity in follicular lymphoma ... by activation-induced deaminase (AID) in *IGH*. Aberrant ... primarily targets noncoding regions causing numerous ... icant “driver” mutations. The quantitative relationship ... SHM, ultradeep sequencing (>20,000-fold coverage) was ... ally targeted by AID (combined 9411 nt), including the 5' ... found in 12/12 FL specimens (median 136 SHMs and 53 ... 1). The number of SNVs at *BCL2* varied widely among ... tential aSHM sites. In contrast, SHM at *IGH* was not

# Expectation: more data → ‘better’ assemblies



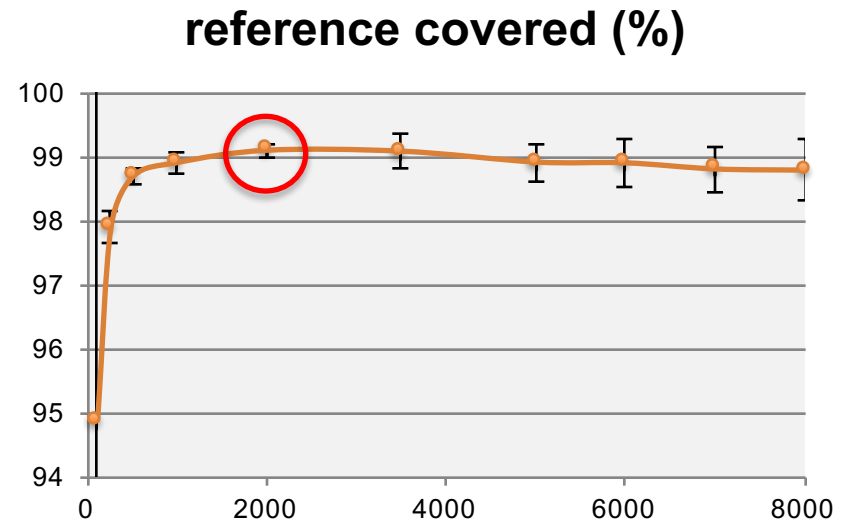
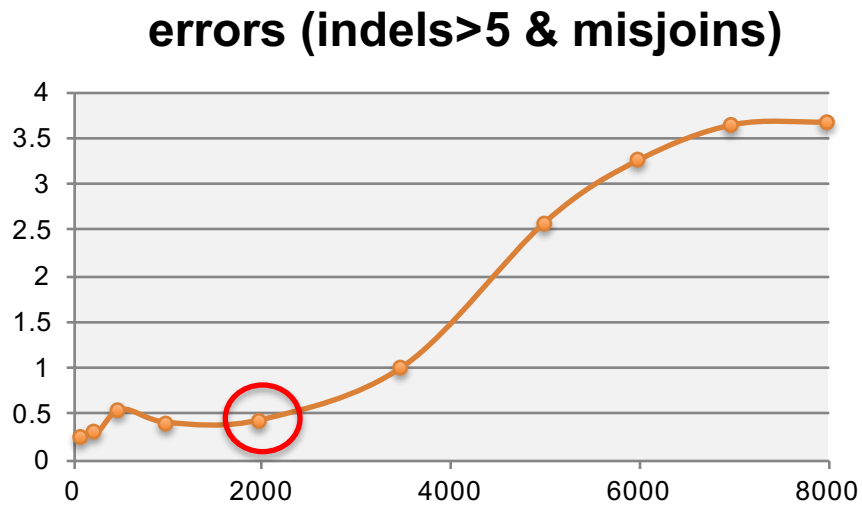
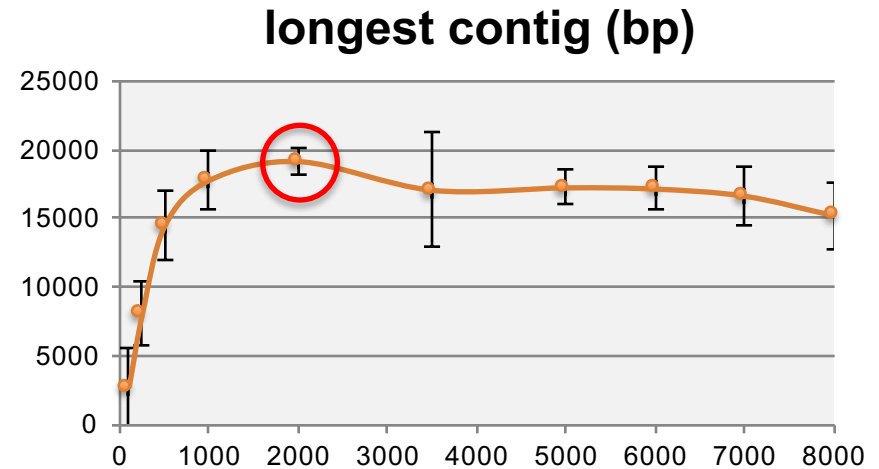
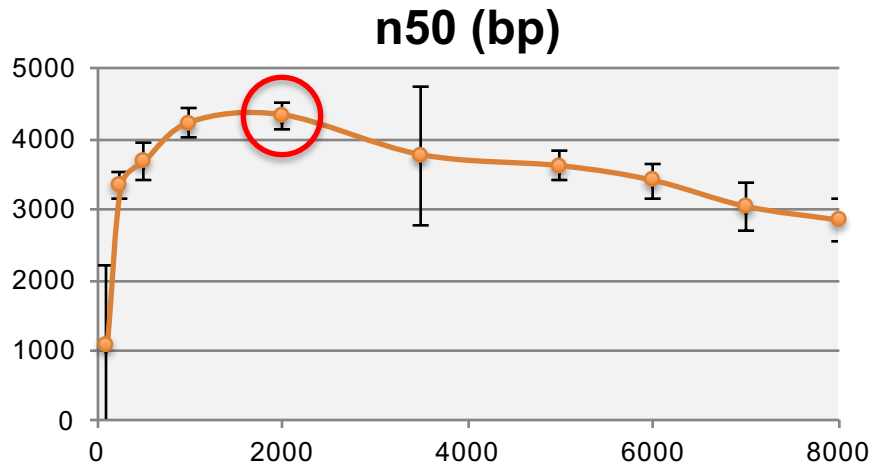
S. Lonardi, H. Mirebrahim, *et al.*, “When less is more: ‘slicing’ sequencing data improves read decoding accuracy and *de novo* assembly quality”, *Bioinformatics*, 2015.

# Expectation: more data → ‘better’ assemblies



S. Lonardi, H. Mirebrahim, *et al.*, “When less is more: ‘slicing’ sequencing data improves read decoding accuracy and *de novo* assembly quality”, *Bioinformatics*, 2015.

# Reality: more data ~~→~~ 'better' assemblies



S. Lonardi, H. Mirebrahim, *et al.*, "When less is more: 'slicing' sequencing data improves read decoding accuracy and *de novo* assembly quality", *Bioinformatics*, 2015.

# More data are not necessarily better: why?

- Possible “suspects”
  - Sequencing errors
  - Highly uneven coverage
  - Read duplication / PCR amplification bias
  - Chimeric reads
  - “Imperfections” in the assembly algorithms

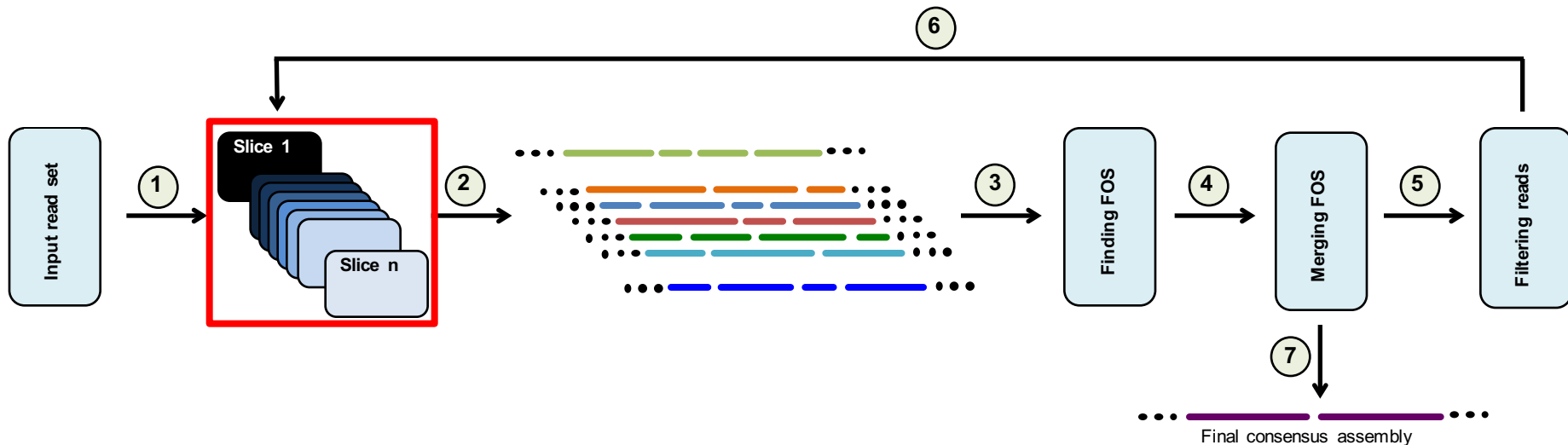
# Possible solutions

- “Classic” error correction
  - based on rare k-mers
  - ineffective for ultra-deep sequencing data
- Down-sampling
  - disregard a fraction of the input reads, according to some predetermined strategy
  - it may remove “critical” reads (i.e., rare error-free reads that can help bridge or fill assembly gaps)
  - not very effective
- SLICEMBLER (next)

# SLICEMBLER algorithm

## 1: “Slice” the input

- The set of input reads is partitioned into  $n$  distinct slices, where  $n = \text{the depth of coverage for the whole input read set} / \text{the desired depth of coverage for each slice}$
- Each slice contains approximately the same number of reads

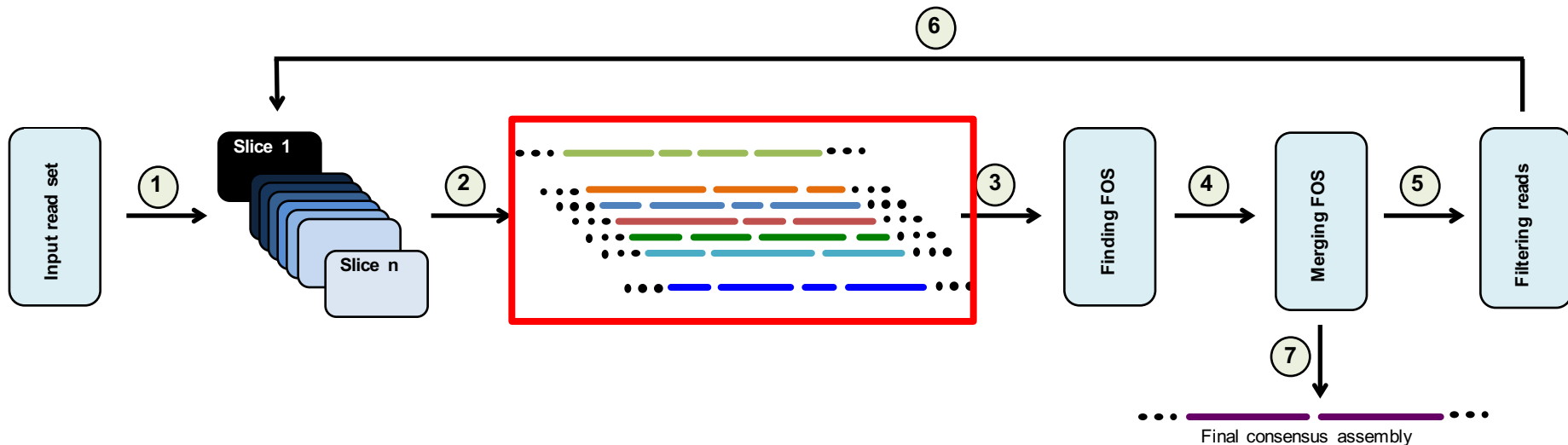




# SLICEMBLER algorithm

## 2: Assemble the reads in each slice

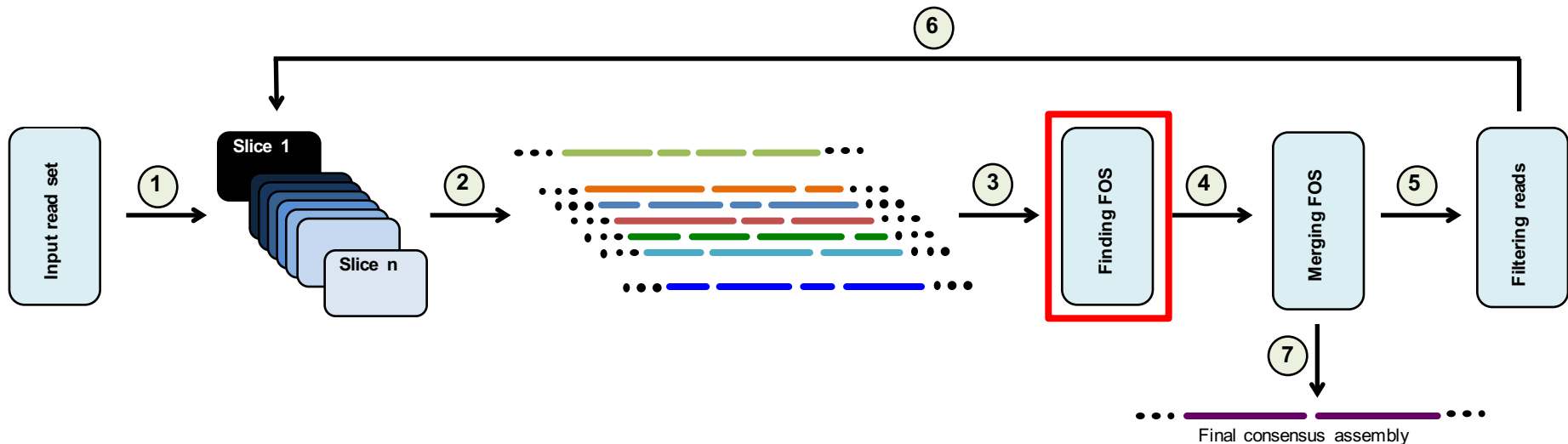
- Each of the  $n$  slices is assembled independently using a standard assembler (Velvet, SPAdes, IDBA, etc.)
- Each assembly is expected to contain a mix of high-quality and low-quality contigs



# SLICEMBLER algorithm

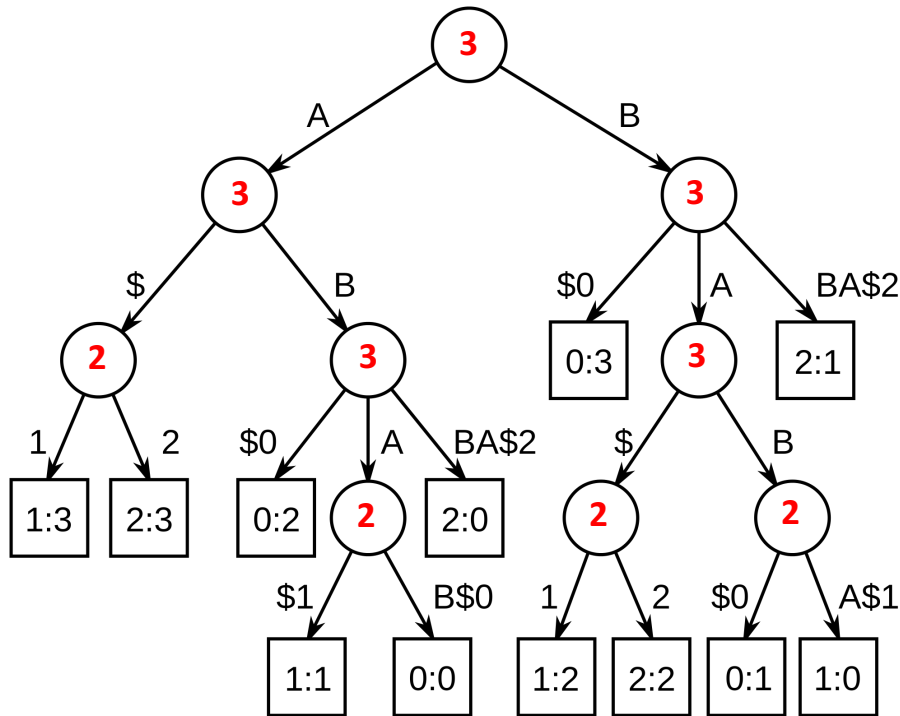
## 3: Find frequently occurring substrings (FOS)

- Identify high-quality contigs (or fraction thereof) or FOS
- Use a generalized suffix-tree for efficiency
- Remove tandem repeats at the end of FOS



# Finding FOS efficiently

**Definition:** Given integers  $k$  and  $l$ , a *FOS* is a maximal substring  $r$  such that  $|r| \geq l$  and it appears in at least  $k$  assemblies

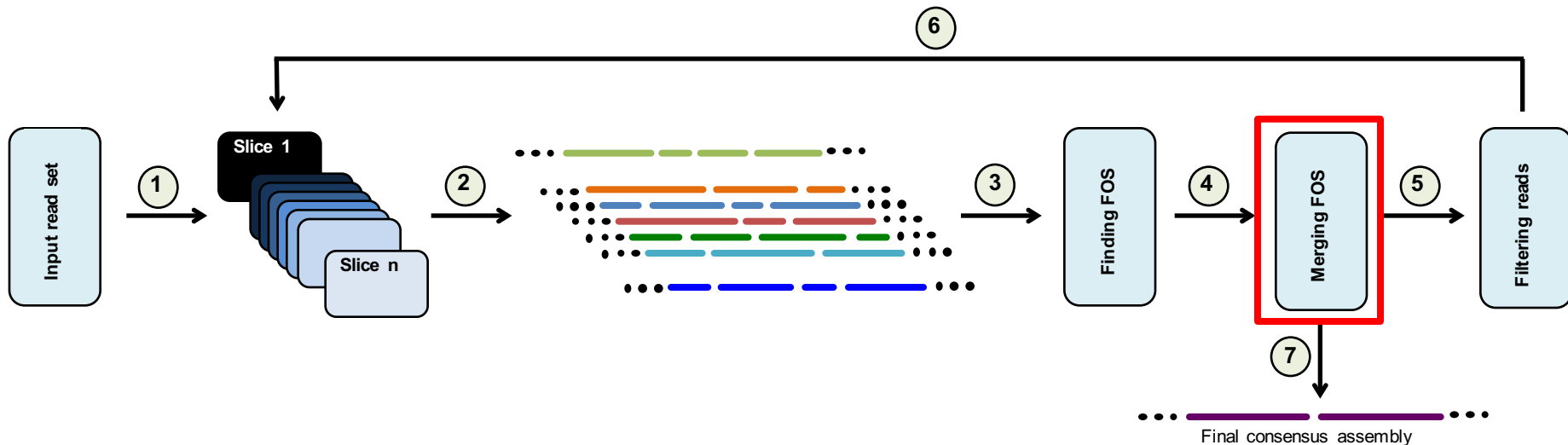


- Build a *generalized suffix tree* on the contigs of the  $n$  assemblies (and their reverse complement)
- Each input assembly is assigned a distinct “color” (Hui, CPM’92)
- Annotate each internal node  $u$  with the number of distinct colors in the subtree rooted at  $u$
- In order to find FOS, determine all the deepest internal nodes (deeper than  $l$ ) which have a color count of at least  $k$
- Building and annotating the suffix tree can be done in linear time

# SLICEMBLER algorithm

## 4: Merge FOS

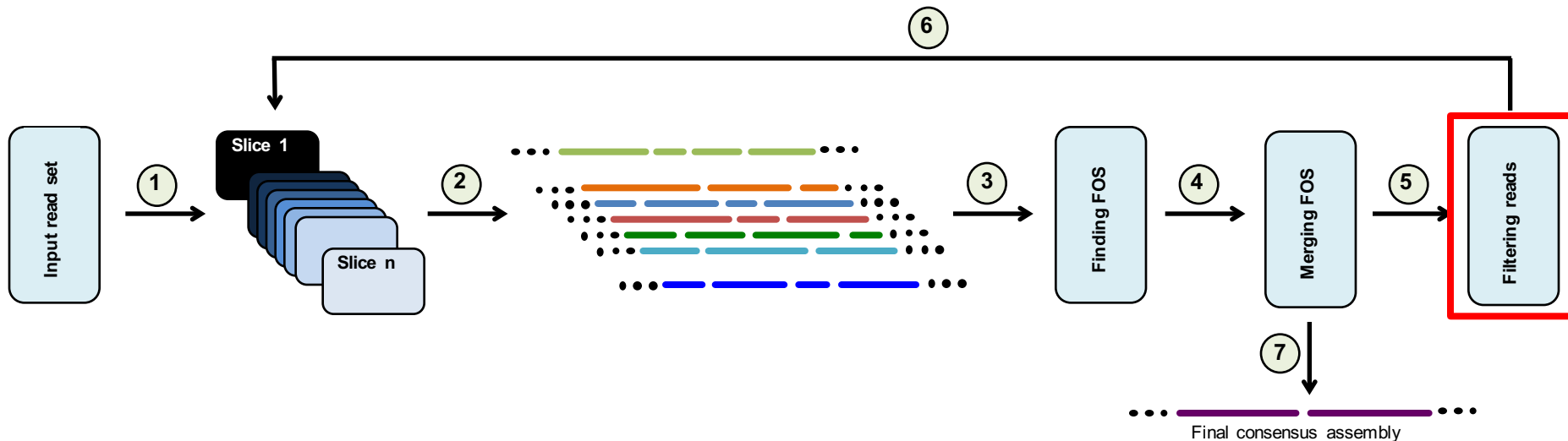
- When detected FOS are overlapping they can be merged to obtain longer FOS
- Merge based on exact suffix-prefix overlap and bridge reads (similar to scaffolding)



# SLICEMBLER algorithm

## 5: Filter reads

- Input reads are mapped to FOS (e.g., BWA)
- Any read that maps to a contig in the current assembly is removed from input (unless it maps close to the end)
- Only the remaining reads are re-assembled in the next iteration



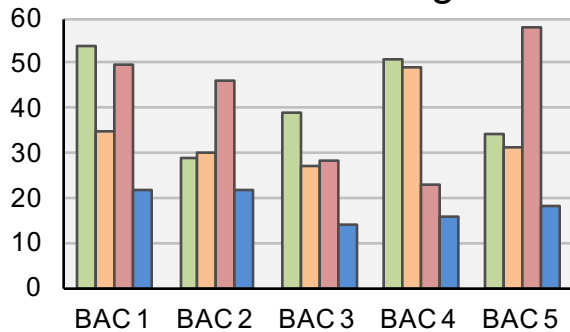
# Experimental results

- Real ultra-deep sequencing data
  - Sequenced 16 barley BACs with Illumina HiSeq
  - Depth of coverage: 8,000x-15,000x
  - Paired-end reads (avg length ~88bp after trimming)
  - Selected 8,000x paired-end reads
  - High-quality references are available for five BACs
- Synthetic ultra-deep reads (*wgsim*)
  - Generated from the reference barley BACs
  - Paired-end reads (2x100 bp)
  - Various levels of coverage and error rates

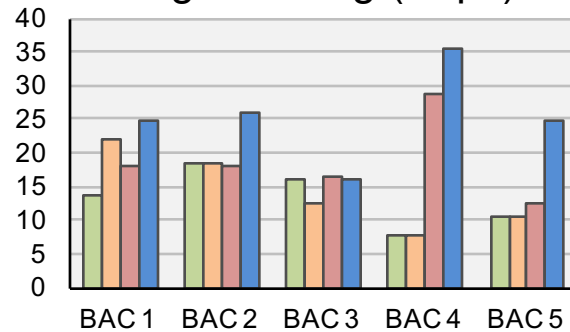
# Experimental results

## (real barley BACs)

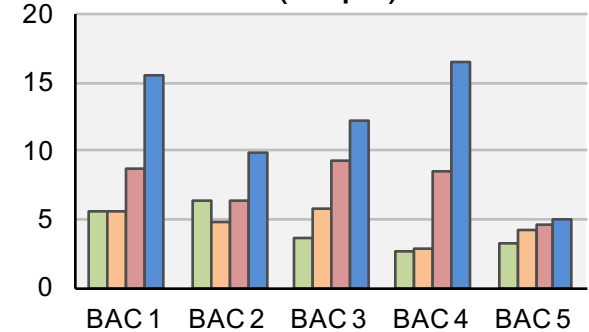
Number of contigs



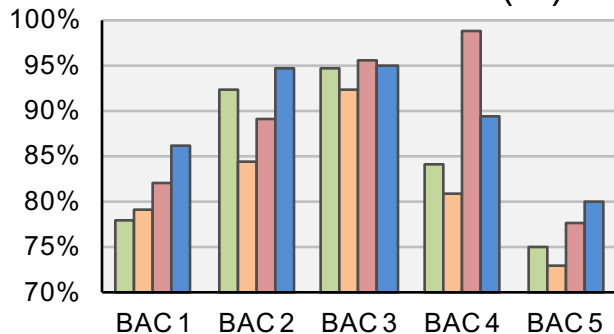
Longest contig (Kbps)



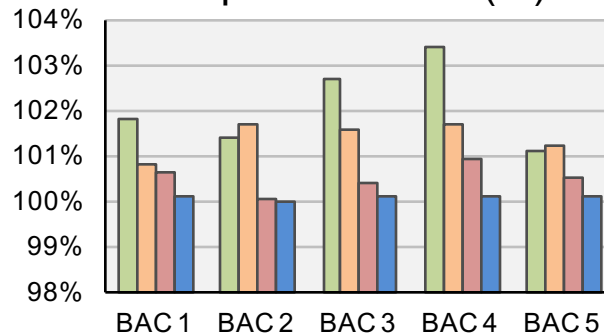
n50 (Kbps)



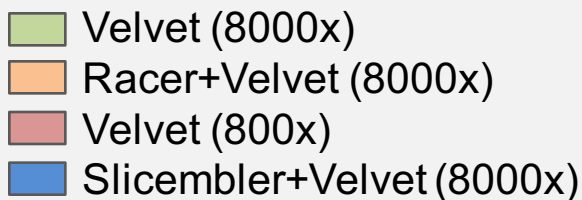
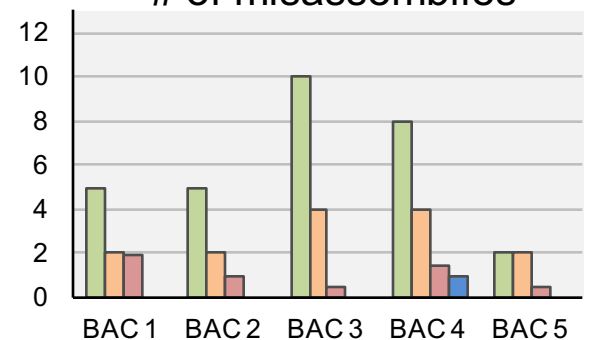
Reference covered (%)



Duplication ratio (%)



# of misassemblies



Statistics collected with QUAST for contigs longer than 500 bp

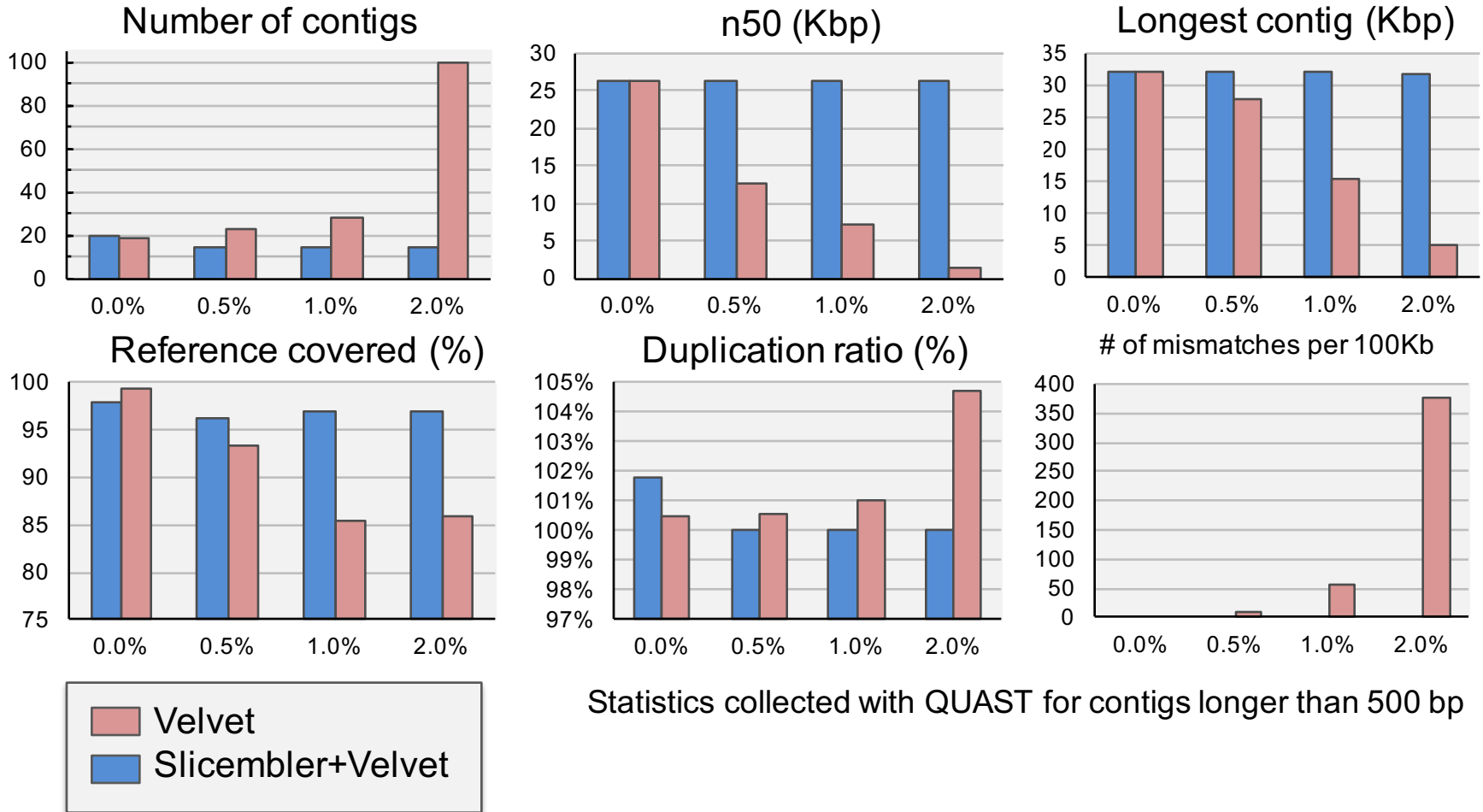
# Assembly quality vs. base assembler (real barley BAC)

<b>Method</b>	<b>Number of contigs</b>	<b>% ref covered</b>	<b>Duplication ratio</b>	<b>Mismatches per 100Kbp</b>	<b>N50</b>	<b>Longest contig</b>
IDBA (8,000x)	34	97.0%	<b>1.010</b>	<b>0.93</b>	7,335	13,889
SLICEMBLER + IDBA (10 slices of 800x)	<b>13</b>	<b>97.0%</b>	<b>1.010</b>	1.1	<b>16,121</b>	<b>31,161</b>
Velvet (8,000x)	39	94.7%	1.027	20.0	3,649	16,048
SLICEMBLER + Velvet (10 slices of 800x)	<b>14</b>	<b>95.1%</b>	<b>1.001</b>	<b>0</b>	<b>12,178</b>	<b>16,128</b>
SPAdes (8,000x)	49	95.7%	<b>1.006</b>	<b>0.94</b>	9,129	21,872
SLICEMBLER + SPAdes (10 slices of 800x)	<b>11</b>	<b>96.9%</b>	1.024	1.2	<b>27,685</b>	<b>31,158</b>
Ray (8,000x)	35	80.0%	1.003	<b>0</b>	3,996	7,186
SLICEMBLER + Ray (10 slices of 800x)	<b>24</b>	<b>88.0%</b>	<b>1.000</b>	<b>0</b>	<b>7,192</b>	<b>12,842</b>

Statistics collected with QUAST for contigs longer than 500 bp



# Assembly quality vs. sequencing error rate (simulated barley BACs)



# Conclusions

- Modern *de novo* genome assemblers seem unable to take advantage of ultra-deep coverage
- SLICEMBLER is an iterative meta-assembler that takes advantage of the whole dataset and due to its “majority voting” scheme
  - Is more resilient to sequencing errors than its base assemblers
  - Almost never incorporates misassemblies in the consensus assembly
- SLICEMBLER is available at [www.slicembler.cs.ucr.edu](http://www.slicembler.cs.ucr.edu)
- SLICEMBLER is slow, but a C++ implementation called SLICEMBLER++, will be available soon

# Thank you



DBI-1062301 and IIS-1302134



NIFA 2009-65300-05645



**USAID**  
FROM THE AMERICAN PEOPLE

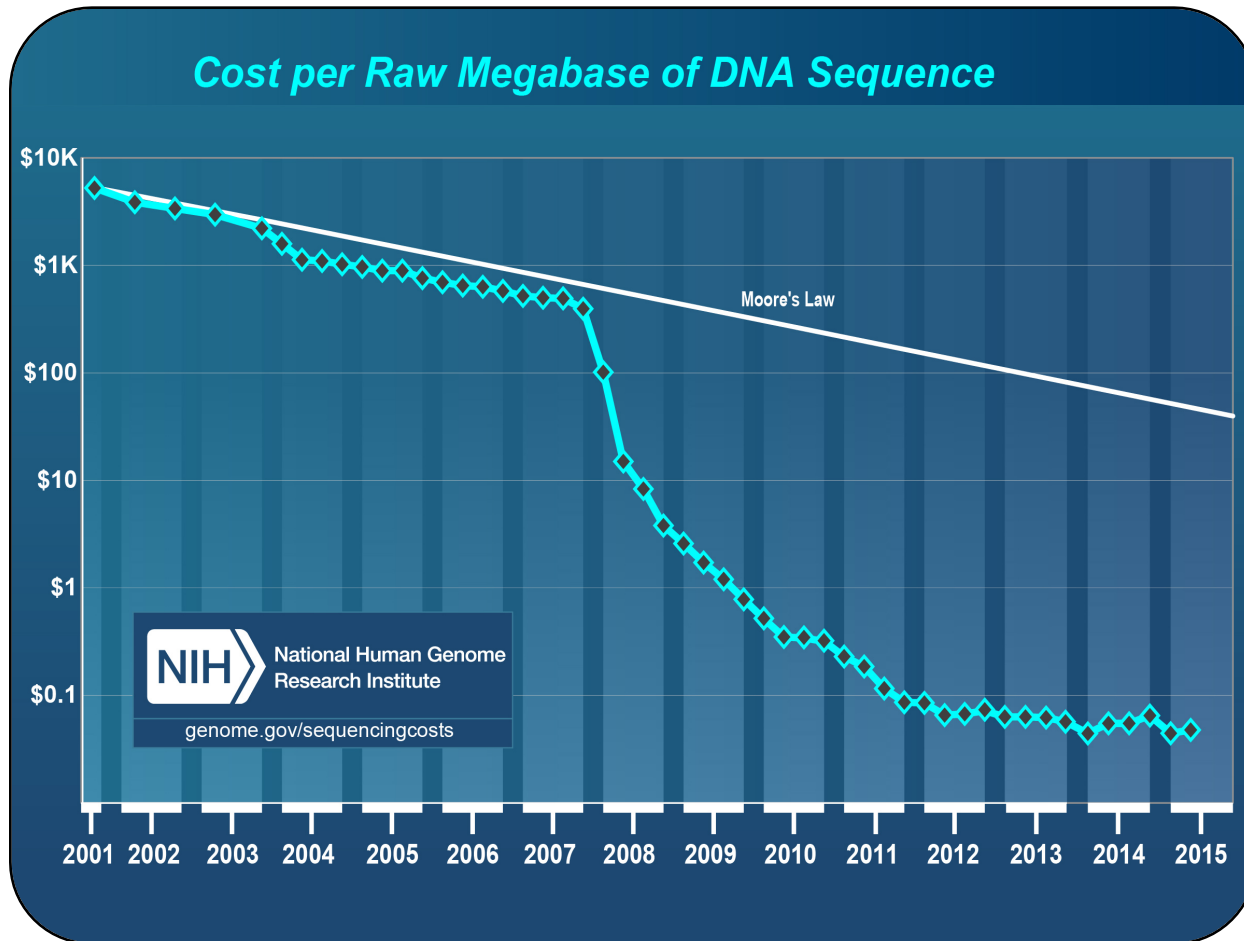
AID-OAA-A-13-00070

[www.slicemblember.cs.ucr.edu](http://www.slicemblember.cs.ucr.edu)



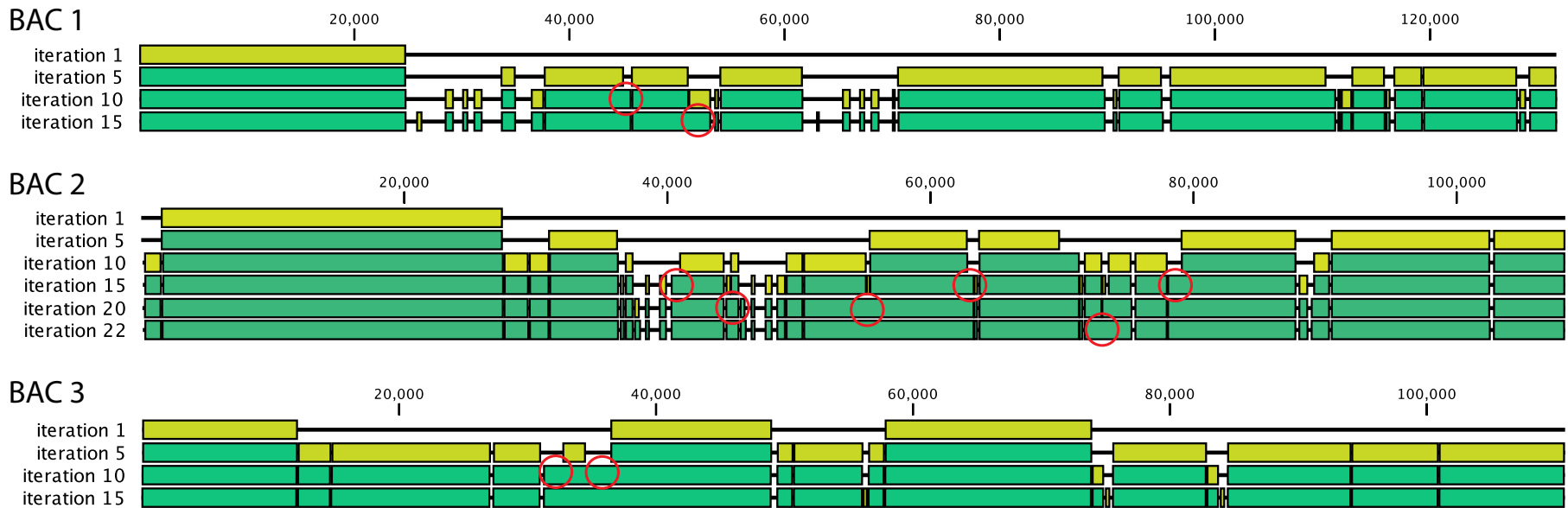
# Will sequencing cost continue to decrease?

*"It's difficult to make predictions, especially about the future"*

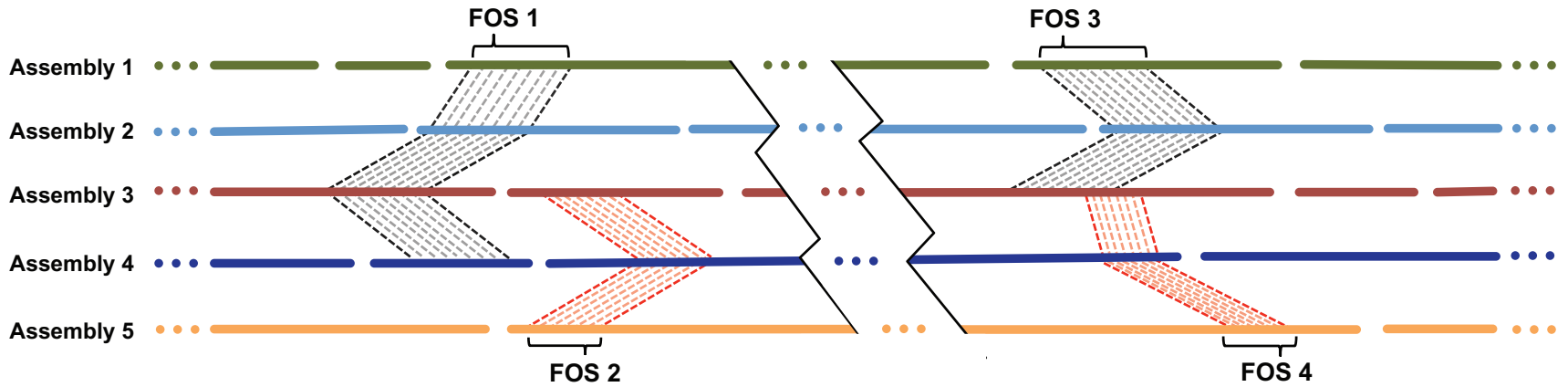


2004 2005 2003 2004 2002 2008 2001 2008 2008 2010 2014 2015 2013 2014 2012

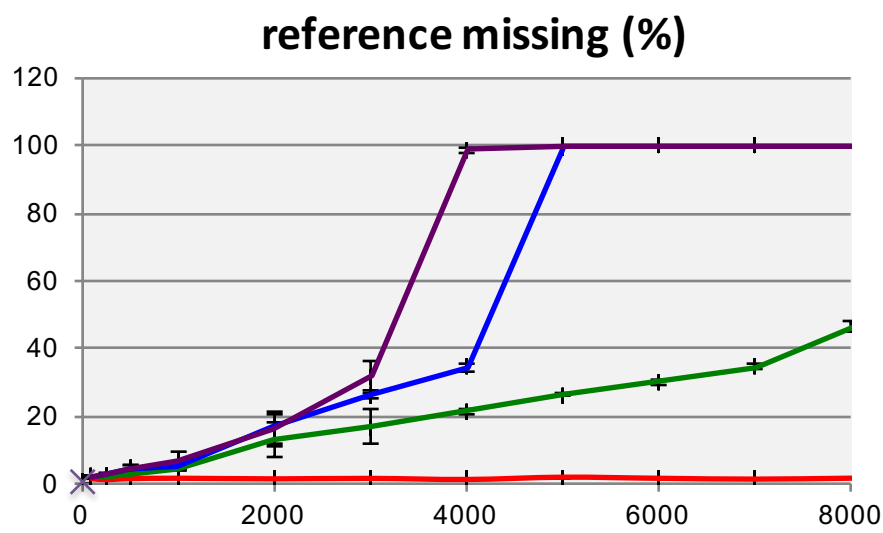
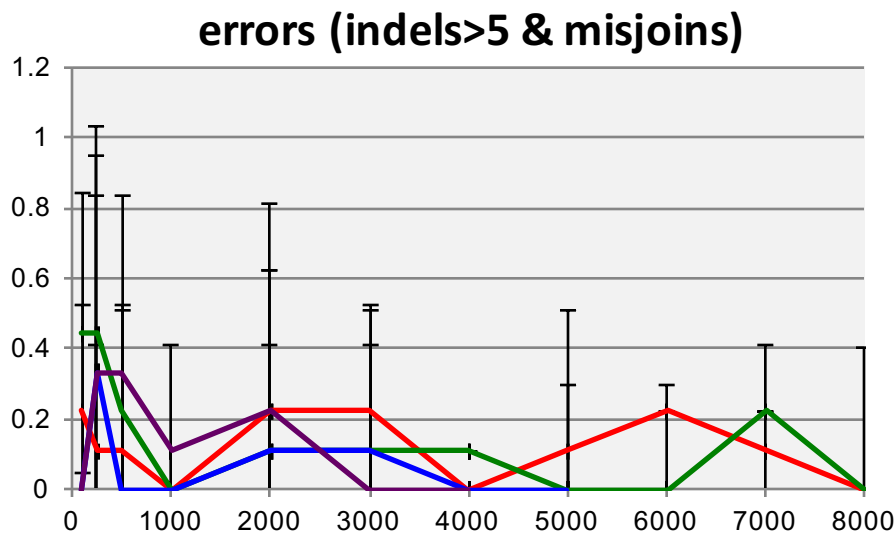
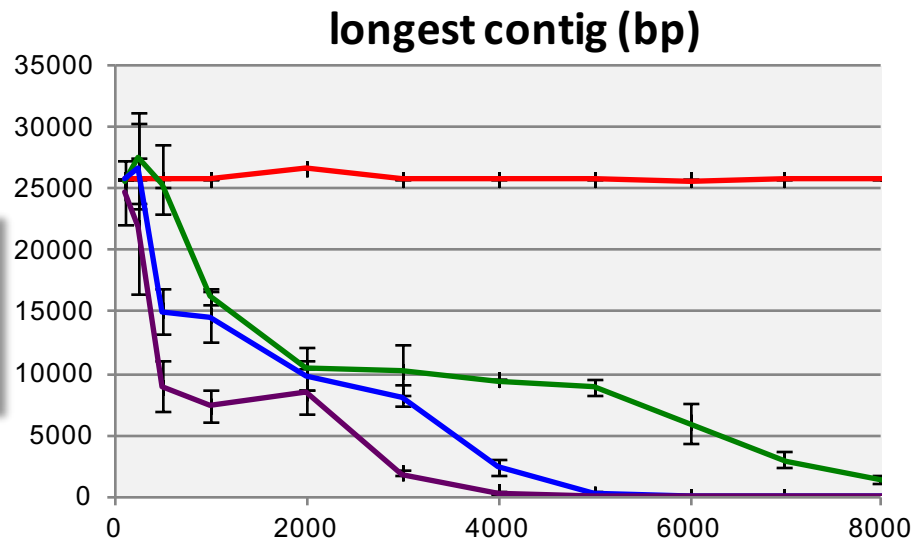
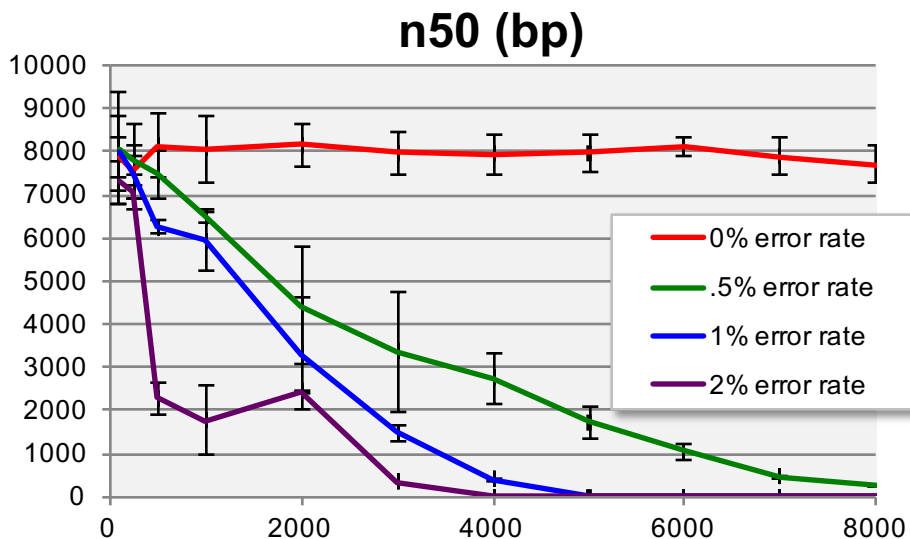
# Experimental results (real barley BACs)



# Frequently occurring substrings (FOS)



# Varying sequencing error rate (Velvet)

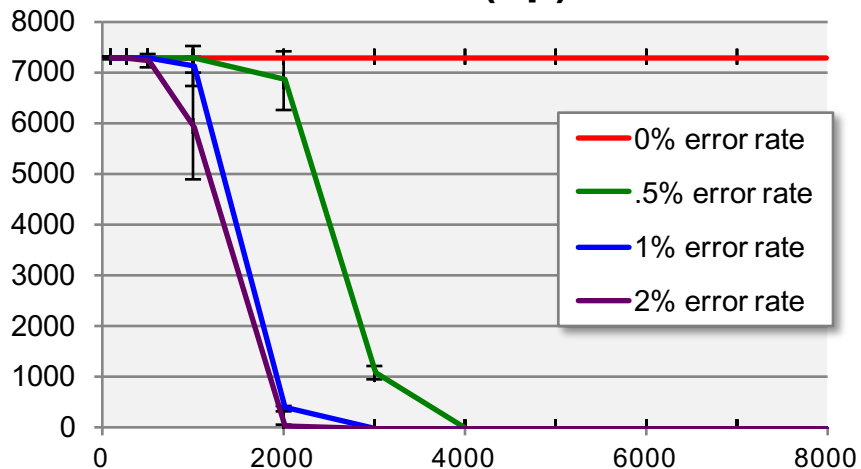


S. Lonardi, H. Mirebrahim, *et al.*, "When less is more: 'slicing' sequencing data improves read decoding accuracy and *de novo* assembly quality", *Bioinformatics*, 2015. <sup>24</sup>

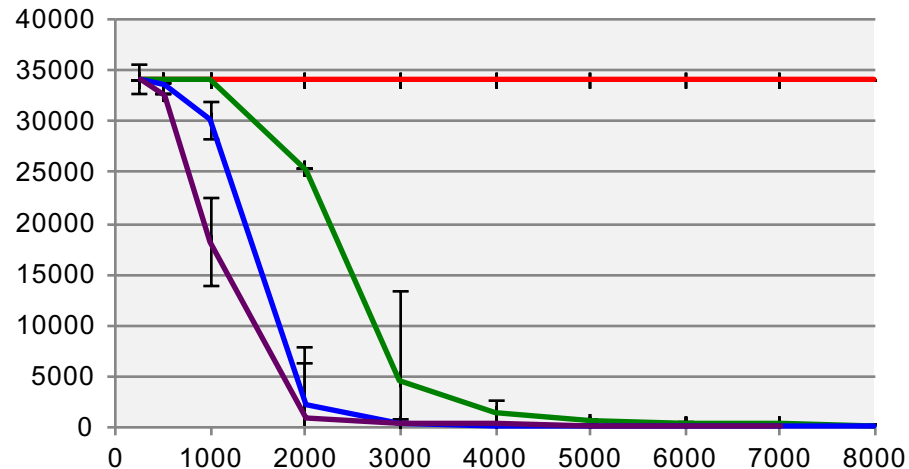


# Varying sequencing error rate (IDBA)

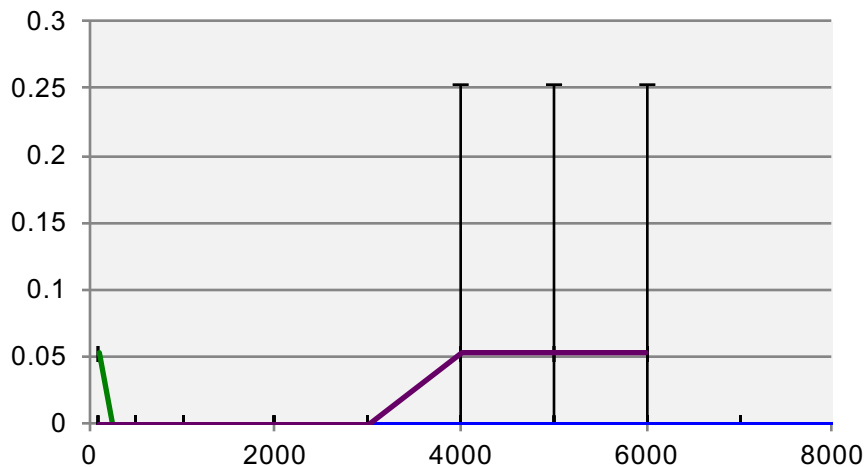
## n50 (bp)



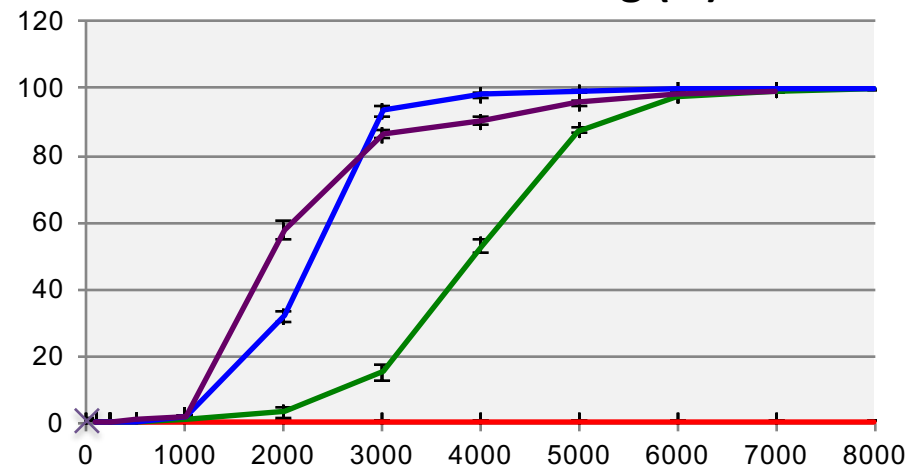
## longest contig (bp)



## errors (indels>5 & misjoins)



## reference missing (%)

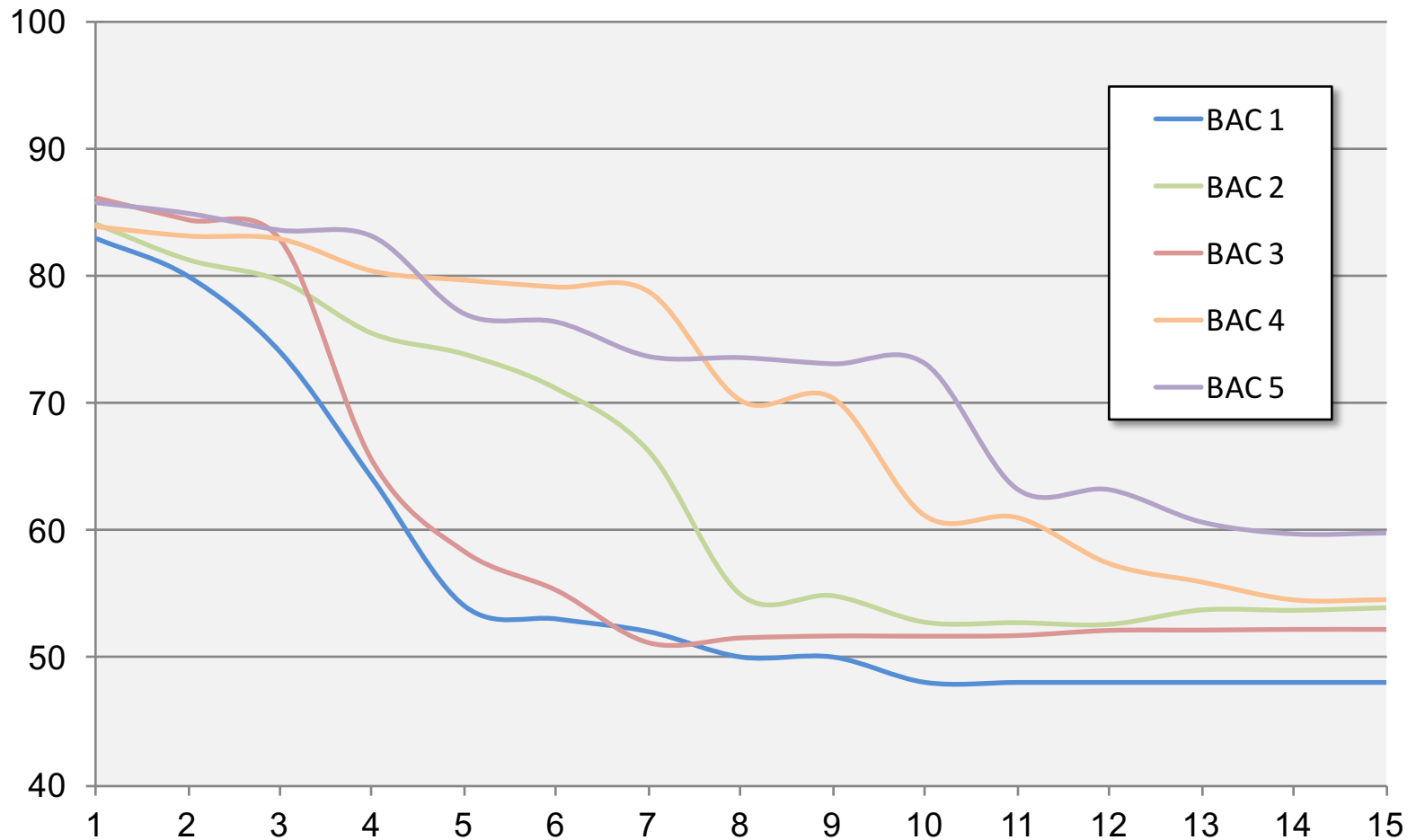


S. Lonardi, H. Mirebrahim, *et al.*, "When less is more: 'slicing' sequencing data improves read decoding accuracy and *de novo* assembly quality", *Bioinformatics*, 2015. <sup>25</sup>

# Sketch of the algorithm

- **partition** the input read set to  $n$  slices, each of which has  $D_s$  coverage
- **while** the cumulative length of FOS is less than the length of genome
  - **assemble** the reads in all  $n$  slices individually
  - **create** a suffix-tree from the  $n$  assemblies and their reverse compl
  - **assign**  $k = n, l = l_{target}/5$
  - **while** ( $l > l_{min}$ )
    - **find** FOS longer than  $l$ , appearing in at least  $k$  assemblies
    - **if** FOS found **then merge** them with the previous FOS, **break**
    - **else if**  $k > n/2$  **then assign**  $k = k-1$ 
      - **else assign**  $l = l/2, k = n$
  - **if** ( $l \leq l_{min}$ ) **and** (no FOS were found) **then break**
  - **map** reads to FOS and **eliminate** mapped reads from the input
- **report** FOS

# Percentage of reads that map exactly to the reference, iteration by iteration



# Assembly quality vs. slice coverage

	500x	1,000x	2,500x	5,000x	7,500x	10,000x
<b>Number of contigs</b>	20	12	11	<b>10</b>	18	38
<b>Longest contig</b>	27,364	31,823	31,946	<b>31,950</b>	21,865	9,425
<b>N50</b>	6,707	26,275	<b>26,288</b>	26,267	12,428	3,643
<b>Percent Refer. Covered</b>	90.6%	88.7%	<b>94%</b>	93.9%	92.9%	84.7%
<b>Duplication ratio</b>	1	1	1	1	1	1
<b>Mismatches per 100kbp</b>	0	0	0	0	0	0

# Another application for SLICEMBLER: co-assembly of single cell sequencing data

