

# Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation

Sergey Koren<sup>1,5</sup>, Brian P. Walenz<sup>1,5</sup>, Konstantin Berlin<sup>2</sup>, Jason R. Miller<sup>3</sup>, Nicholas H. Bergman<sup>4</sup> and Adam M. Phillippy<sup>1</sup>

<sup>1</sup>Genome Informatics Section, Computational and Statistical Genomics Branch, National Human Genome Research Institute, National Institutes of Health, Bethesda, Maryland 20892, USA;

<sup>2</sup>Invincea Incorporated, Fairfax, Virginia 22030, USA;

<sup>3</sup>J. Craig Venter Institute, Rockville, Maryland 20850, USA;

<sup>4</sup>National Biodefense Analysis and Countermeasures Center, Frederick, Maryland 21702, USA

Presented By  
**Anmol Sharma<sup>1</sup>**

<sup>1</sup>MSc Student in Computing Science  
Medical Image Analysis Laboratory  
Simon Fraser University  
Burnaby, BC Canada.

20 September 2018

# Outline

## 1 Problem Statement

- Overview

## 2 Background

- DNA Sequencing
- Sequencing Technologies

## 3 Canu Pipeline

- Overview
- Correction Stage
- Trimming Stage
- Assembly Stage

## 4 Results

- Performance and Comparison

## 5 Conclusion

# Outline

## 1 Problem Statement

### ■ Overview

## 2 Background

### ■ DNA Sequencing

### ■ Sequencing Technologies

## 3 Canu Pipeline

### ■ Overview

### ■ Correction Stage

### ■ Trimming Stage

### ■ Assembly Stage

## 4 Results

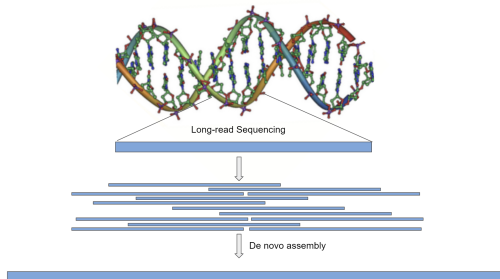
### ■ Performance and Comparison

## 5 Conclusion

# Problem Statement

## Overview

- *De novo* genome assembly from long-read, single molecule sequencing.
- Due to low accuracy of long-read sequencing technologies, efficient and accurate assembly of large repeats and closely related haplotypes remains challenging.
- **Question** : How do we accurately assemble the whole genome using noisy, inaccurate long-reads ?



# Outline

- 1 Problem Statement
  - Overview
- 2 Background
  - DNA Sequencing
  - Sequencing Technologies
- 3 Canu Pipeline
  - Overview
  - Correction Stage
  - Trimming Stage
  - Assembly Stage
- 4 Results
  - Performance and Comparison
- 5 Conclusion

# Background

## DNA Sequencing

- DNA is the basic building block of any organism which contains genetic “code” which makes all chemical activity in an organism possible.
- A code is analogous to a long sentence, and “bases” are the words that form that sentence. “Sequencing” a DNA refers to mapping out the exact sequence of bases as they are laid out in the DNA.
- Mapping the bases accurately is a required precursor for any downstream analysis such as in forensics, medical science and mutations.
- In recent years, advances in sequencing technology (like PacBio and Oxford Nanopore) have enabled us to sequence the DNA faster, accurately, and with relatively low cost.

# Outline

- 1 Problem Statement
  - Overview
- 2 Background
  - DNA Sequencing
  - Sequencing Technologies
- 3 Canu Pipeline
  - Overview
  - Correction Stage
  - Trimming Stage
  - Assembly Stage
- 4 Results
  - Performance and Comparison
- 5 Conclusion

# Background

## Sequencing Technologies

- **PacBio** from Pacific Biosciences : The Sequel System is based on Single Molecule, Real-Time (SMRT).
  - Capable of reading between 10-20kbp per read.
  - Workflow of <1 day.
- **Oxford Nanopore** : Direct, real time sequencing technology.
  - Capable of reading upto 2mbp per read.
  - Available in handheld versions as well as industrial sized versions.



# Background

## Outstanding Issues

- Despite rapid advances in sequencing technology there are still outstanding issues.
- Average number of bases in the DNA varies tremendously between different organisms, ranging from 4 million (E. coli K12) to 3 billion (human genome).
- The current sequencing technologies are only capable of randomly reading small subsequences of DNA from random locations, that too with errors.
- This is akin to a large jigsaw puzzle, with multiple overlapping pieces which may not even be correct, and no reference image to facilitate assembly of jigsaw pieces.

# Outline

## 1 Problem Statement

### ■ Overview

## 2 Background

### ■ DNA Sequencing

### ■ Sequencing Technologies

## 3 Canu Pipeline

### ■ Overview

### ■ Correction Stage

### ■ Trimming Stage

### ■ Assembly Stage

## 4 Results

### ■ Performance and Comparison

## 5 Conclusion

# Canu Pipeline

## Overview of the System

- Canu pipeline is modular assembly of three stages :
  - Correction
  - Trimming
  - Assembly
- Each stage can be run independently or in sequence.
- Canu supports both single computer node and compute clusters, with support for most major job schedulers (SLURM, LSF, SGE).
- Summary statistics are output as the job progresses.
- Each stage uses it's own indexed database, which disregards the previous input data once created.

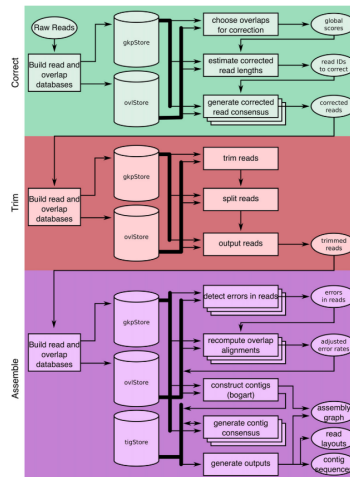


FIGURE – Canu whole pipeline

# Outline

- 1 Problem Statement
  - Overview
- 2 Background
  - DNA Sequencing
  - Sequencing Technologies
- 3 Canu Pipeline
  - Overview
  - **Correction Stage**
  - Trimming Stage
  - Assembly Stage
- 4 Results
  - Performance and Comparison
- 5 Conclusion

# Correction Stage

## Overview

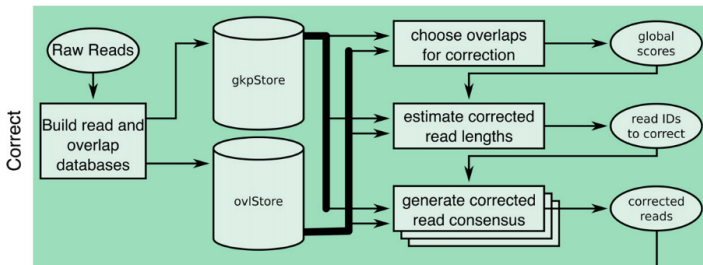


FIGURE – Correction Stage

# Correction Stage

Identify Overlaps : MinHash Alignment Process (MHAP)

- MinHash Alignment Process (MHAP) was proposed in (Berlin 2015<sup>1</sup>) for aligning overlapping noisy, long reads using probabilistic, locality-sensitive hashing.
- MHAP integrated with Celera Assembler was able to perform *reference grade* de novo assemblies of many genome sequences.
- It uses the idea of MinHash Sketches to perform rapid overlapping of noisy reads.

---

1. Berlin, Konstantin, et al. "Assembling large genomes with single-molecule sequencing and locality-sensitive hashing." *Nature biotechnology* 33.6 (2015) : 623.

# Correction Stage

Identify Overlaps : MHAP Example

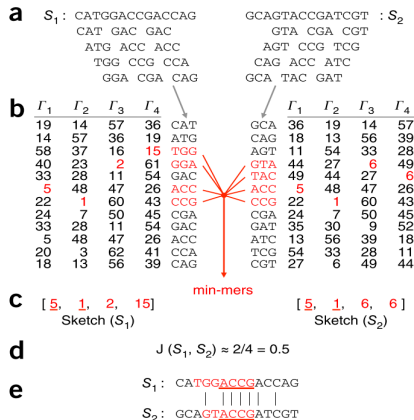
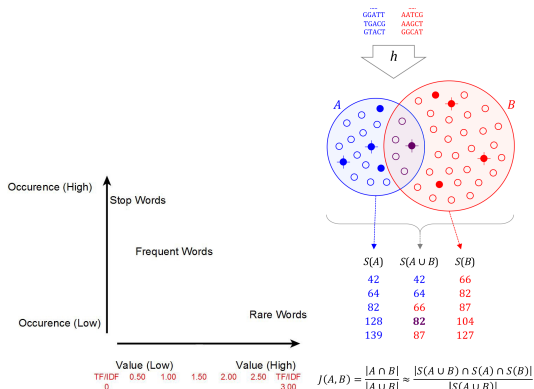


FIGURE – MHAP Example

# Correction Stage

Identify Overlaps : MinHash Alignment Process (MHAP) with *tf-idf*

- *Canu* uses a modified version of vanilla MHAP.
- In *Canu*, MHAP uses a two-stage overlap filter.
  - STAGE 1 uses *tf-idf* weighting to prefer informative, non-repetitive *k*-mers which increases sensitivity to true overlaps and reduces number of false, repetitive overlaps considered.
  - STAGE 2 MHAP uses "bottom-sketch" strategy similar to the one proposed in (Ondov 2016)<sup>2</sup>. Used to check the quality of overlap.



2. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, Phillippy AM. 2016. Mash : fast genome and metagenome distance estimation using MinHash. Genome Biol 17 : 132



# Correction Stage

## Read Correction

- Long read sequencing exhibit high error rate (5-15% depending on the technology).
- There is a need to to correct each of the individual reads using a consensus based approach.
- *Canu* uses all-vs-all overlap information to correct individual reads.
- It uses *two* filtering steps to determine which overlaps should be chosen to correct each individual read.
  - STAGE 1 (Global Filter) Each read chooses to provide correction evidence for *C* other reads.
  - STAGE 2 (Local Filter) Each read accepts/rejects the correction evidence supplied by other reads.
- This strategy was first proposed by (Koren et al. 2012)<sup>3</sup> as PBcR (PacBio Corrected Reads Pipeline), which performs hierarchical correction and assembly of single-molecule reads.
- Corrected reads are then generated using “falcon\_sense” algorithm, originally proposed in (Chin et al. 2016)<sup>4</sup>

3. Koren, Sergey, et al. "Hybrid error correction and de novo assembly of single-molecule sequencing reads." *Nature biotechnology* 30.7 (2012) : 693.

4. Chin CS, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, Dunn C, O'Malley R, Figueroa-Balderas R, Morales-Cruz A, et al. 2016. Phased diploid genome assembly with single-molecule real-time sequencing. *Nat Methods* 13 : 1050–1054.

# Outline

- 1 Problem Statement
  - Overview
- 2 Background
  - DNA Sequencing
  - Sequencing Technologies
- 3 Canu Pipeline
  - Overview
  - Correction Stage
  - **Trimming Stage**
  - Assembly Stage
- 4 Results
  - Performance and Comparison
- 5 Conclusion

# Trimming Stage

## Overview

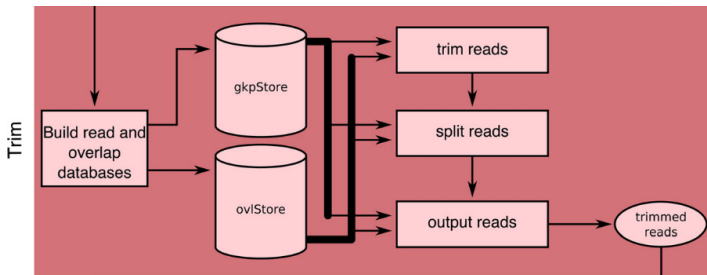


FIGURE – Trim Stage

# Trimming Stage

## Overlap Based Trimming (OBT)

- Once reads are corrected in the Correction stage, and new database of reads and overlaps is created once again in Trimming stage through overlapInCore, the reads are trimmed.
- Trimming is performed according to the algorithm first described in (Miller et al. 2008)<sup>6</sup>.
- 5 ■ The reads are trimmed to the largest portion covered to at least depth  $C$  by overlaps of at most  $E$  error, and minimum length  $L$ .
- Parameters are technology dependent, and are chosen empirically.
- A second pass is made to detect hairpin adapters and chimeras in the reads. Once detected, the reads are trimmed to largest supported region.

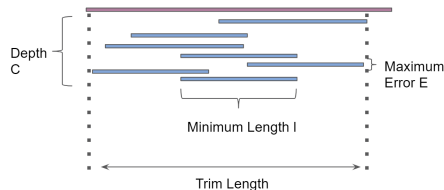


FIGURE – First pass to trim reads

5. Miller JR, Delcher AL, Koren S, Venter E, Walenz BP, Brownley A, Johnson J, Li K, Mobarry C, Sutton G. 2008. Aggressive assembly of pyrosequencing reads with mates. Bioinformatics 24 : 2818–2824.

# Outline

- 1 Problem Statement
  - Overview
- 2 Background
  - DNA Sequencing
  - Sequencing Technologies
- 3 Canu Pipeline
  - Overview
  - Correction Stage
  - Trimming Stage
  - **Assembly Stage**
- 4 Results
  - Performance and Comparison
- 5 Conclusion

# Assembly Stage

## Overview

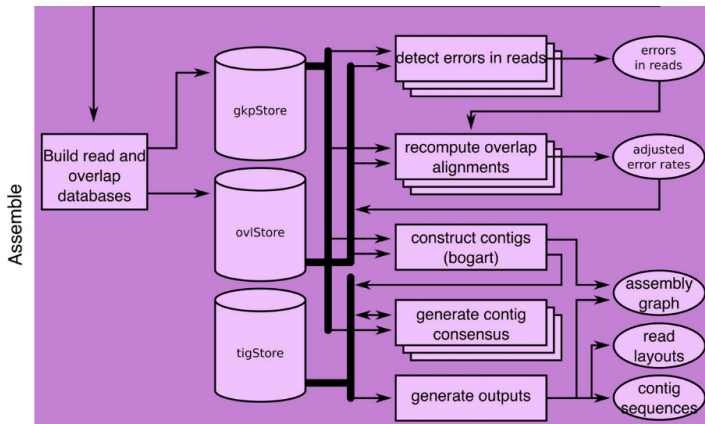


FIGURE – Assembly Stage

# Assembly Stage

## Overlap Error Adjustment

- After trimming and before graph construction, *Canu* recomputes overlaps and makes final attempt at detecting sequencing errors.
- The intuition is to improve separation between true sequencing differences (repeats and haplotypes) and false differences due to random error.
- Each read is corrected by the majority vote of it's overlapping reads.
- However the reads are **not changed**, since the overlaps have already been constructed.
- Reported error rates for each overlap is adjusted **had the changes in reads were made**.
- Two passes are made through the reads. The first pass finds the errors in reads, and the second pass temporarily applies the fixes in order to calculate reported error rates.
- This algorithm was first used in (Holt et al. 2002)<sup>6</sup>.

6. Holt RA, Subramanian GM, Halpern A, Sutton GG, Charlab R, Nusskern DR, Wincker P, Clark AG, Ribeiro JC, Wides R, et al. 2002. The genome sequence of the malaria mosquito *Anopheles gambiae*. Science 298 : 129–149

# Assembly Stage

## Graph Construction

- Graph construction uses the “best overlap graph” (BOG) strategy from (Miller et al. 2008)<sup>a</sup>.
  - Overlaps are defined as “containment”.
    - If all bases in one read are aligned to another read, or *dovetail*, if involving only the ends of both reads.
  - A “best overlap” is the longest dovetail overlap to a given read end.
- 7
- In original paper, all overlaps are picked up to a user specified cut-off threshold.
  - In “BOGART” (Canu’s version of “best overlap graph”), the best overlaps are chosen after several filtering steps.
    - This removes high-error overlaps, potential chimeric reads, and reads whose overlaps indicate a possible sequence anomaly.
  - BOGART results in a cleaner and more accurate graph construction.

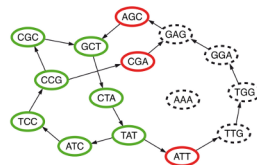


FIGURE — Overlap graph

7. Miller JR, Delcher AL, Koren S, Venter E, Walenz BP, Brownley A, Johnson J, Li K, Mobarry C, Sutton G. 2008. Aggressive assembly of pyrosequencing reads with mates. *Bioinformatics* 24 : 2818–2824



# Assembly Stage

## Contig Consensus

- Contig consensus is performed using a slightly modified version of “pbdagcon” algorithm proposed in (Chin et al. 2013)<sup>8</sup>.
- A “template sequence” is constructed for each contig by splicing reads together from approximate positions based on BOG.
- The template is accurate within individual reads, but may still contain indel errors at read boundaries.
- To correct this, all reads in contig are aligned to template sequence using Myers’ O(ND) algorithm<sup>9</sup>, and added to a DAG.
- The DAG is used to call a consensus sequence as described in (Chin et al. 2013)<sup>10</sup>.

8. Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, Clum A, Copeland A, Huddleston J, Eichler EE, et al. 2013. Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. Nat Methods 10 : 563–569

9. Myers EW. 1986. An O(ND) difference algorithm and its variations. Algorithmica 1 : 251–266

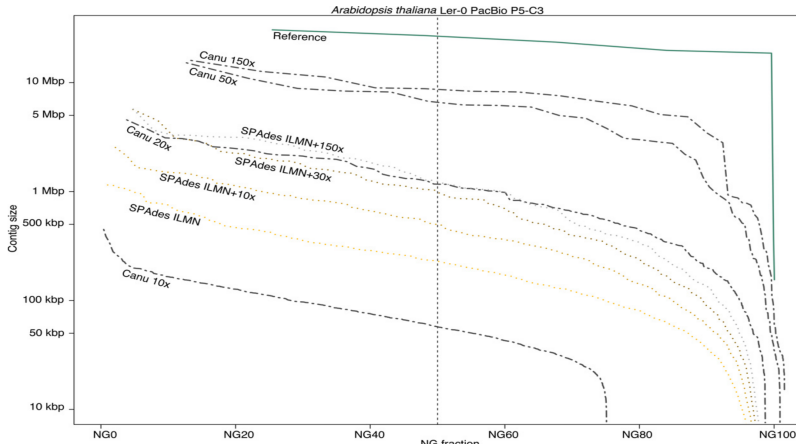
10. Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, Clum A, Copeland A, Huddleston J, Eichler EE, et al. 2013. Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. Nat Methods 10 : 563–569

# Outline

- 1 Problem Statement
  - Overview
- 2 Background
  - DNA Sequencing
  - Sequencing Technologies
- 3 Canu Pipeline
  - Overview
  - Correction Stage
  - Trimming Stage
  - Assembly Stage
- 4 Results
  - Performance and Comparison
- 5 Conclusion

# Results

Performance of Canu as a function of coverage



# Results

## Canu on PacBio Data

**Table 1.** Canu is fastest for generating a high-quality polished assembly from PacBio data

Genome	Asm/Polish	Max (Mbp)	NG50 (Mbp)	% Ref	No. of breakpoints	Time (CPU h)	% Idy
<i>Escherichia coli</i>	Canu+Quiver	<b>4.68</b>	<b>4.68</b>	<b>100%</b>	<b>0</b>	12.25	<b>99.9999%</b>
	FALCON+Quiver	4.64	4.64	<b>100%</b>	2	25.14	99.9998%
	Miniasm+Quiver	4.64	4.64	99.99%	2	31.93	99.9998%
	SPAdes	4.64	4.64	<b>100%</b>	<b>0</b>	<b>4.09</b>	99.9972%
<i>Drosophila melanogaster</i>	Canu+Quiver	<b>25.78</b>	<b>21.31</b>	<b>97.47%</b>	1025	<b>1396.52</b>	99.9795%
	FALCON+Quiver	23.08	9.84	96.12%	1054	2305.92	<b>99.9813%</b>
	Miniasm+Quiver	15.85	5.84	96.51%	<b>752</b>	1484.33	<b>99.9813%</b>
	Canu+Quiver	<b>15.95</b>	<b>8.31</b>	<b>82.94%</b>	220	<b>925.31</b>	<b>99.0710%</b>
<i>Arabidopsis thaliana</i>	FALCON+Quiver	15.94	8.17	82.72%	222	1132.25	<b>99.0710%</b>
	Miniasm+Quiver	11.61	5.07	82.88%	<b>205</b>	976.43	<b>99.0710%</b>
	Canu+Quiver	5.34	2.35	<b>99.70%</b>	139	410.07	<b>99.9745%</b>
<i>Caenorhabditis elegans</i>	FALCON+Quiver	4.99	1.88	98.82%	<b>138</b>	<b>397.40</b>	99.9735%
	Miniasm+Quiver	<b>5.85</b>	<b>2.96</b>	99.44%	141	526.16	99.9706%
	Canu+Quiver	80.08	<b>21.95</b>	<b>86.84%</b>	1105	<b>22,749.71</b>	99.8081%
CHM1	FALCON+Quiver	52.34	9.46	86.58%	<b>1082</b>	68,789.00	<b>99.8086%</b>

# Results

## Canu on Oxford Nanopore Data

**Table 2.** Canu consistently assembles complete genomes from only Oxford Nanopore data

Genome	Asm/Polish	No. of contigs	Max (Mbp)	% Ref	No. of breakpoints	Time (CPU h)	% Idy
<i>Escherichia. coli</i> MAP005	Canu+Nanopolish	<b>(1)</b>	<b>4.64</b>	<b>99.98%</b>	2	376.87	<b>99.43%</b>
	FALCON+Nanopolish	105	0.42	23%	2	<b>106.2</b>	99.41%
	Miniasm+Nanopolish	3	3.40	99.96%	<b>0</b>	2344.02	99.36%
<i>E. coli</i> MAP006-1	Canu+Nanopolish	<b>(1)</b>	4.63	99.80%	<b>0</b>	<b>167.04</b>	<b>99.81%</b>
	FALCON+Nanopolish	<b>(1)</b>	4.63	99.86%	<b>0</b>	207.45	99.78%
	Miniasm+Nanopolish	<b>(1)</b>	<b>4.66</b>	<b>99.97%</b>	2	1801.02	99.72%
<i>E. coli</i> MAP006-2	Canu+Nanopolish	<b>(1)</b>	4.64	99.91%	<b>2</b>	<b>168.69</b>	<b>99.78%</b>
	FALCON+Nanopolish	<b>(1)</b>	4.64	<b>99.94%</b>	<b>2</b>	196.16	99.76%
	Miniasm+Nanopolish	<b>(1)</b>	<b>4.65</b>	99.70%	4	1482.95	99.69%
<i>E. coli</i> MAP006-PCR-1	Canu+Nanopolish	<b>(1)</b>	<b>4.64</b>	99.95%	<b>0</b>	<b>164.08</b>	<b>99.84%</b>
	FALCON+Nanopolish	<b>(1)</b>	4.63	99.80%	2	168.37	99.82%
	Miniasm+Nanopolish	3	2.15	<b>99.96%</b>	<b>0</b>	1338.28	99.77%
<i>E. coli</i> MAP006-PCR-2	Canu+Nanopolish	<b>(1)</b>	4.64	99.99%	2	<b>206.09</b>	<b>99.85%</b>
	FALCON+Nanopolish	<b>(1)</b>	4.64	<b>100.00%</b>	2	212.89	99.84%
	Miniasm+Nanopolish	<b>(1)</b>	<b>4.65</b>	99.98%	<b>0</b>	1669.83	99.81%
<i>Bacillus anthracis</i>	Canu+Nanopolish	<b>(2)</b>	5.20	<b>99.77%</b>	<b>0</b>	894.40	99.14%
	FALCON+Nanopolish	31	0.47	86.29%	<b>0</b>	<b>795.93</b>	<b>99.17%</b>
	Miniasm+Nanopolish	4	<b>5.22</b>	97.21%	<b>0</b>	5094.90	99.05%
<i>Yersinia pestis</i>	Canu+Nanopolish	<b>(4)</b>	4.67	<b>99.97%</b>	<b>11</b>	<b>254.25</b>	<b>99.76%</b>
	FALCON+Nanopolish	<b>(4)</b>	<b>4.68</b>	<b>99.97%</b>	12	295.01	99.72%
	Miniasm+Nanopolish	9	2.69	99.91%	<b>11</b>	2000.16	99.65%

# Results

## Canu vs Hybrid Methods

**Table 3.** Nanopore assemblies exceed hybrid methods in continuity and match their quality when polished with Illumina data

Genome	Asm/Polish	No. of contigs	Max (Mbp)	% Ref	No. of breakpoints	Time (CPU h)	% Idy
<i>E. coli</i> MAP005	Canu+Pilon	<b>(1)</b>	<b>4.65</b>	99.99%	2	10.98	99.9873%
	FALCON+Pilon	105	0.42	23.04%	2	4.36	99.9550%
	Miniasm+Pilon	3	3.40	90.62%	42	<b>3.15</b>	97.3878%
	SPAdes	<b>(1)</b>	4.64	<b>100.00%</b>	<b>0</b>	3.61	<b>99.9989%</b>
<i>E. coli</i> MAP006-1	Canu+Pilon	<b>(1)</b>	4.63	99.82%	<b>0</b>	5.89	<b>99.9995%</b>
	FALCON+Pilon	<b>(1)</b>	4.63	99.86%	<b>0</b>	7.3	99.9964%
	Miniasm+Pilon	<b>(1)</b>	<b>4.66</b>	96.97%	21	<b>3.14</b>	99.6118%
	SPAdes	<b>(1)</b>	4.64	<b>100.00%</b>	<b>0</b>	3.65	99.9965%
<i>E. coli</i> MAP006-2	Canu+Pilon	<b>(1)</b>	<b>4.64</b>	99.94%	2	3.92	<b>99.9987%</b>
	FALCON+Pilon	<b>(1)</b>	<b>4.64</b>	99.94%	2	3.93	99.9933%
	Miniasm+Pilon	<b>(1)</b>	<b>4.64</b>	97.98%	26	<b>2.73</b>	99.6336%
	SPAdes	<b>(1)</b>	<b>4.64</b>	<b>100.00%</b>	<b>0</b>	3.56	99.9965%
<i>E. coli</i> MAP006-PCR-1	Canu+Pilon	<b>(1)</b>	<b>4.64</b>	99.95%	<b>0</b>	4.15	<b>99.9993%</b>
	FALCON+Pilon	<b>(1)</b>	4.63	99.80%	2	3.55	99.9969%
	Miniasm+Pilon	3	2.16	98.41%	12	<b>2.15</b>	99.6734%
	SPAdes	2	3.95	<b>100.00%</b>	<b>0</b>	3.56	99.9965%
<i>E. coli</i> MAP006-PCR-2	Canu+Pilon	<b>(1)</b>	4.64	<b>100.00%</b>	2	6.16	<b>99.9992%</b>
	FALCON+Pilon	<b>(1)</b>	4.64	<b>100.00%</b>	2	9.22	99.9963%
	Miniasm+Pilon	<b>(1)</b>	<b>4.65</b>	98.57%	20	<b>2.69</b>	99.6734%
	SPAdes	<b>(1)</b>	4.64	<b>100.00%</b>	<b>0</b>	4.00	99.9965%
<i>B. anthracis</i>	Canu+Pilon	<b>(2)</b>	5.21	99.77%	1	65.01	99.8476%
	FALCON+Pilon	31	0.48	86.31%	<b>0</b>	14.95	99.8888%
	Miniasm+Pilon	4	<b>5.25</b>	79.36%	44	<b>4.9</b>	92.2732%
	SPAdes	6	4.13	<b>100.00%</b>	<b>0</b>	8.47	<b>99.9948%</b>
<i>Y. pestis</i>	Canu+Pilon	<b>(4)</b>	<b>4.66</b>	<b>99.83%</b>	23	17.92	99.8946%
	FALCON+Pilon	<b>(4)</b>	4.64	99.65%	26	10.63	99.8715%
	Miniasm+Pilon	9	2.70	93.76%	42	<b>8.68</b>	98.7866%
	SPAdes	29	0.37	95.99%	<b>15</b>	17.08	<b>99.9559%</b>

# Conclusion

Canu

- Canu is a three stage sequence assembly system that brings together improved state-of-art methods.
- Canu assembles high-quality polished assembly from PacBio data, and highly continuous genomes from Oxford Nanopore data.
- When compared to methods which perform the same task, Canu retrieved the highest coverage of the reference sequence, and also dramatically reduced the computational time (CPU hours).
- Canu provides a robust and efficient method for sequence assembly from long read sequence data, and enables high quality downstream analysis of the sequences due to its redundant error correction steps that lead to accurate and complete assemblies of the genome.

Thank You !



Questions ?