

# Genome Assembly Tutorial

Deb Triant

22 October 2018

Programming for Biology  
Cold Spring Harbor, NY

# Genome Assembly Workshop



Genome assembly with PacBio data  
using Canu

Why are we using this program?

- \*Relatively\* user-friendly and easy to install
- Not computationally intensive with small data sets
- Well maintained and documented

# Genome assemblers

- Other useful assemblers:

• **Illumina data:**

• **w2rap-contigger**

<https://github.com/bioinfologics/w2rap-contigger>

- can take lots of computer power to run
- works with single paired-end library

• **soapdenovo2**

<https://sourceforge.net/projects/soapdenovo2>

- relatively easy to install and run
- works with large genomes

# Genome assemblers

## PacBio data:

falcon & falcon-unzip

<https://github.com/PacificBiosciences/FALCON>

<http://profs.scienze.univr.it/delledonne/Papers/2016%20Chin%20NMethods.pdf>

Very powerful assemblers but not easy to install or use.

Quiver/Arrow/pbalign - genome alignment and polishing. Part of PacBio Genomic Consensus package.

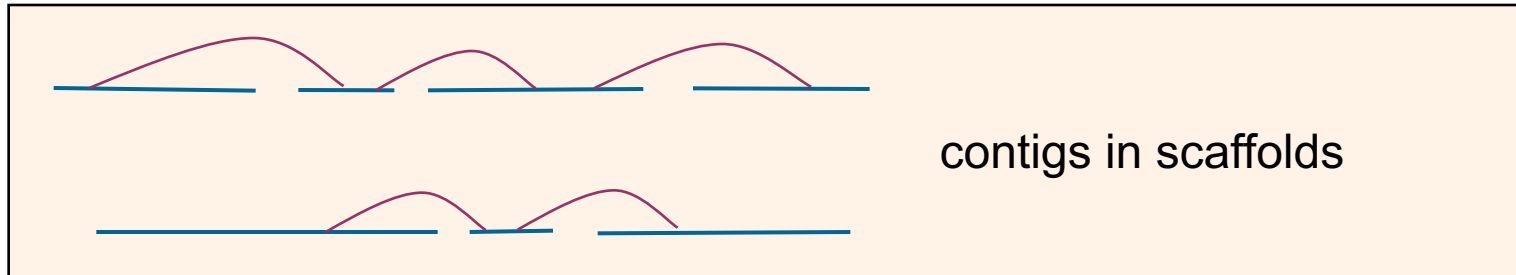
<https://github.com/PacificBiosciences/GenomicConsensus>

Again, powerful but not the easiest to use.

Canu

<http://genome.cshlp.org/content/27/5/722>

# Linear assemblies



**contig**: a contiguous sequence of bases....

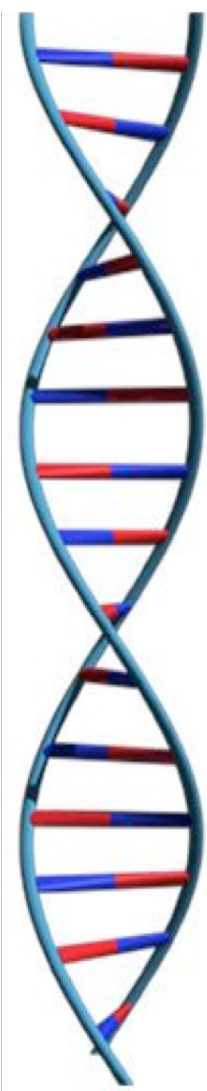
```
CTGCCCCCTGTGCCAATGGGTTTGAGGCTCTTCCCACCTTCCTTTTCTATTAGATTCAATGTATCTGGTTTTATGTTGAGG
TCCTAGATCCACTTGGACTTGAGCTTTGTACAAGATGACATATATAGGTCTGTTTTTATTCTTCTACATACAGACAGCCA
GTTATAACCAGCACCATTTATTGAAGACACTTTCTTTATTCCATTGTATATTTTTTTACTTCCTTGTCAAAAATCAAGTGA
CCATGAGTATGTGGTTTCATTTCTGGGTCTTCAATTGTATTCCATTAGTCAACATATCTGTCTCTGTACCAATACCATGC
```

**scaffold**: a sequence of contigs, separated by gaps....

```
TCCTAGATCCACTTGGACTTGAGCTTTGTACAAGATGACATATATAGGTCTGTTTTTATTCTTCTACATACAGACAGCCA
GTTATAACCAGCACCATTTATTGAAGACACTTTCTTTATTCCATTGTATATTTTTTTACTTCCTTGTCAAAAATCAAGTGA
CCATGAGTATGTGGTTTCATTTCTGGGTCTTCAATTGTATTCCATTAGTCAACATATCTGTCTCTGTACCAATACCATGC
NNNNNNNNAGTTTTTACCACAATTGCTCTATAGTAAAGCTTGAGGTCAGGGTTGGTGATCCCTCCAGCCATTCTTTTCATT
ATTAAGAATTGTTTTCCCTAGTCTGGGTTTTTTTGCTTTTCCAGGCGAATTTGAGAATTGCTCTTTCCATGTCTTTGAAGA
ATTGTGTTNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNGGGATTTTGATGGGGTTTGCATTGAATCTGTAGATT
GTCTTTGGTAAGATGGTTAGTTTTACTATGTAAATTCTGCCAATCCACAAGCATGGGAGCGCTCTCCATTTTCTGAGATC
TTCTTCAATTTCTTTCTTGAGAACTTGAAGTTATTGTCATACA
```

Number of Ns = predicted gap size, with error bars (can't be displayed in fasta format)

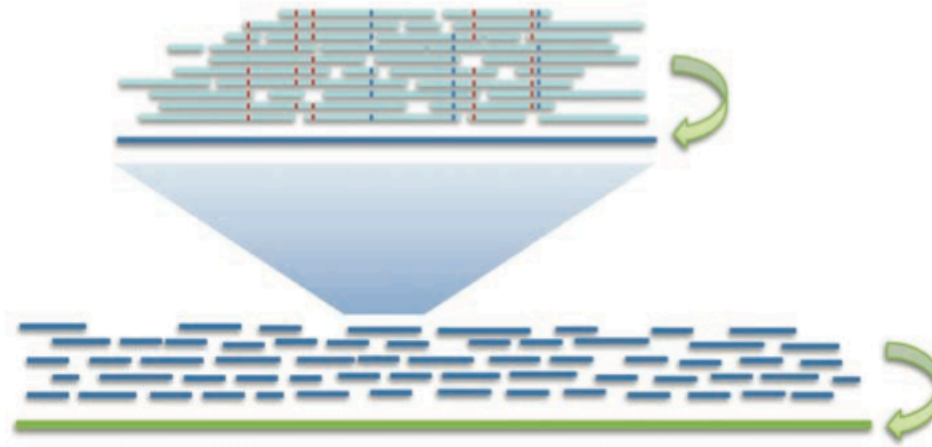
# Long read sequencing



Longest sequencing  
reads as seed data set

Use seed dataset  
to map shorter  
reads

Perform preassembly:  
construct preassembled  
reads from seed reads  
through consensus  
procedure



Preassembled reads  
used for genome  
assembly

Error correction: mapping high-quality  
short PacBio reads to long reads

**PB-only Correction & Polishing**

Chin et al (2013) *Nature Methods*. 10:563–569

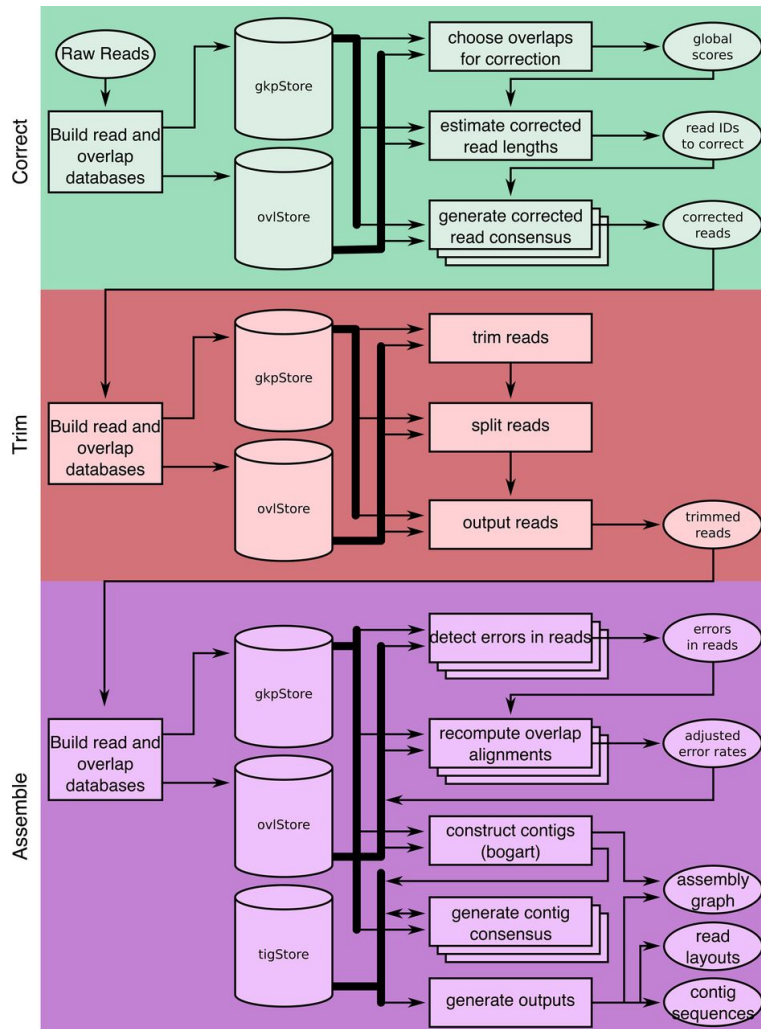
Canu: “unitigs”

# Using Canu

- Canu is an assembler that specializes in PacBio or Oxford Nanopore long-read data
- Derived from Celera assembler
- Corrects reads, trims suspect regions (e.g., adaptors) then assembles the corrected reads
- Can be run with one command to do all steps or each step can be run separately (correcting, trimming and assembling).
- Koren S, Walenz BP, Berlin K, Miller JR, Phillippy AM. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Research*  
<http://canu.readthedocs.io/en/latest/index.html>

**A full Canu run includes three stages: correction (green), trimming (red), and assembly (purple).**

Generates k-mer histogram and conducts all-vs-all overlaps



**Correct:**

Select best overlaps to generate corrected reads. Uses longest 40X reads for correction

**Trim:**

Identifies unsupported regions and trims reads to longest supported range

**Assemble:**

Identify sequencing errors, constructs best overlap graph and outputs contigs.



# Using Canu

- Outputs contigs and unitigs.
- The unitig construction task finds sets of overlaps that are consistent, and uses those to place reads into a multialignment layout. The layout is then used to generate a consensus sequence for the unitig or a “high-confidence contig”

# Installing Canu

- We will use the sample Pac Bio *E. coli* data found on the canu homepage: <http://canu.readthedocs.io>  
25x *E.coli* fastq files with quality scores.

## To install Canu:

Login as admin

Make a directory for installation

```
$ git clone https://github.com/marbl/canu.git
```

```
$ cd canu/src
```

```
$ make -j 4 .... Compiler error?
```

# Installing Canu

Compiler error?

```
$ brew install gcc
```

```
$ which gcc
```

Try making again:

```
$ make -j 4
```

Where is it installed?

Executables found in `PATHWAY/canu/Darwin-amd64/bin`

# Installing Canu

Need gnuplot for plotting:

```
$ which gnuplot
```

```
$ brew install gnuplot@4
```

If you need to have gnuplot@4 first in your PATH run:

```
echo 'export  
PATH="/usr/local/opt/gnuplot@4/bin:$PATH" '  
>> ~/.bash_profile
```

# Installing Canu

## Getting Help

`/canu/documentation/source - help docs`

`$ canu command line options`

`$ canu -options parameter options`

# Installing Canu

## Download canu test E.coli dataset:

<http://canu.readthedocs.io/en/latest/quick-start.html?highlight=25x>

- 25X subset (223Mb)

Where do you want them to download?

```
$ curl -L -o pacbio.fastq  
http://gembox.cbcb.umd.edu/mhap/raw/ecoli_p6_  
25x.filtered.fastq
```

- o naming downloaded file `pacbio.fastq`

- L if page has moved to a new location, curl will redo the request at the new site

How many reads within the file? How many lines?

# Running Canu

**screen** - unix command to run programs in multiple windows or “screens”

- If your computer breaks connection, your program will continue run on server

\$ **screen** - to open a new screen

\$ **ctrl-a d** to suspend screen

- Program will continue to run

\$ **screen -r** - to re-open existing screen session

\$ **screen -r #####** - use number to re-open if multiple sessions running

# Running Canu

- We will run default with correcting, trimming and assembling all in one command:

```
$ /your pathway/Deb/canu/Darwin-amd64/bin/canu  
  -p ecoli -d ecoli-pacbio genomeSize=4.8m -pacbio-raw \  
your_path/canu/pacbio.fastq saveReads=true \  
  gnuplot=/usr/local/opt/gnuplot@4/bin/gnuplot  
2>CanuRun_20181022.log **to save output
```

\*where are reads....include pathway

-p assembly-prefix

-d output directory

**input types:**

-pacbio-raw -pacbio-corrected -nanopore-raw -nanopore-corrected fastq  
or fasta format

\***gnuplot option** specific to our exercise

**\*\*USE THE HELP DOCS!!!!**



# Running Canu - Manual assembly

## Correct, Trim and Assemble Manually

- You can also do the three top-level tasks by hand. This would allow trying multiple construction parameters on the same set of corrected and trimmed reads.

Previous command to correct, trim & assemble:

```
canu -p ecoli -d ecoli-full genomeSize=4.8m -pacbio-raw \  
your_path/canu/pacbio.fastq saveReads=true \  
2>full_out_102517.log **to save output
```

# Using Canu – Manual assembly

First correct the raw reads:

- `canu -correct -p ecoli -d ecoli-manual genomeSize=4.8m \`  
`-pacbio-raw your_path/canu/pacbio.fastq \`  
`saveReads=true 2>correct_out_102517.log`

Then trim the output of the correction:

- `canu -trim -p ecoli -d ecoli-manual genomeSize=4.8m \`  
`-pacbio-corrected ecoli-manual/ecoli.correctedReads.fasta.gz \`  
`saveReads=true 2>trim_out_102517.log`

\*corrected reads created during first step

# Using Canu - manual assembly

Finally, assemble the output of the trimming twice using two error rates (can use as many as you like):

- `canu -assemble -p ecoli -d ecoli-erate-0.013 genomeSize=4.8m \`  
`correctedErrorRate=0.013 -pacbio-corrected ecoli- \`  
`manual/ecoli.trimmedReads.fasta.gz saveReads=true \`  
`2>assemble-0.013_out_102517.log`
- `canu -assemble -p ecoli -d ecoli-erate-0.025 genomeSize=4.8m \`  
`correctedErrorRate=0.025 -pacbio-corrected ecoli- \`  
`manual/ecoli.trimmedReads.fasta.gz saveReads=true \`  
`gnuplotImageformat=svg 2>assemble-0.025_out_102517.log`

\*trimmed reads created during second step

\*The error rate specifies the difference in overlap between two corrected reads which is typically <1% (canu default value 0.045) for PacBio data and <2% (canu default 0.144) for Nanopore data (<1% on newest chemistries). Higher rate, more sensitive.

\*Notice in the output there are separate directories for each error rate you specify.

# Running Canu

- We will run default with correcting, trimming and assembling all in one command:

```
$ /your pathway/Deb/canu/Darwin-amd64/bin/canu  
  -p ecoli -d ecoli-pacbio genomeSize=4.8m -pacbio-raw \  
your_path/canu/pacbio.fastq saveReads=true \  
  gnuplot=/usr/local/opt/gnuplot@4/bin/gnuplot  
2>CanuRun_20181022.log **to save output
```

\*where are reads....include pathway

-p assembly-prefix

-d output directory

**input types:**

-pacbio-raw -pacbio-corrected -nanopore-raw -nanopore-corrected fastq  
or fasta format

\***gnuplot option** specific to our exercise

**\*\*USE THE HELP DOCS!!!!**

# Canu output files - all found in helpdocs

- **ecoli\*.report** - assembly analysis log

## READS

- **ecoli\*.correctedReads.fasta.gz** - The sequences after correction, trimmed and split based on consensus evidence. Typically >99% for PacBio and >98% for Nanopore but it can vary based on your input sequencing quality
- **ecoli\*.trimmedReads.fasta.gz** - The sequences after correction and final trimming. The corrected sequences above are overlapped again to identify any missed hairpin adapters or bad sequence that could not be detected in the raw sequences.

## SEQUENCE

- **ecoli\*.contigs.fasta** - Everything which could be assembled and is part of the primary assembly, including both unique and repetitive elements.
- **ecoli\*.unitigs.fasta** - Contigs, split at alternate paths in the graph
- **ecoli\*.unassembled.fasta** - Reads and low-coverage contigs which could not be incorporated into the primary assembly.

**How many contigs to we have? How many unitigs?**

# Canu output files - all found in helpdocs

## SEQUENCE

- **ecoli\*.contigs.fasta** - Everything which could be assembled and is part of the primary assembly, including both unique and repetitive elements.
- **ecoli\*.unitigs.fasta** - Contigs, split at alternate paths in the graph
- **ecoli\*.unassembled.fasta** - Reads and low-coverage contigs which could not be incorporated into the primary assembly.

The header line for each sequence provides some metadata on the sequence:

```
>tig00000001 len=576432 reads=573 covStat=249.02 gappedBases=no class=contig  
suggestRepeat=no suggestCircular=no
```

```
tig # len=<integer> reads=<integer> covStat=<float> gappedBases=<yes|no>  
class=<contig|bubble|unassm> suggestRepeat=<yes|no> suggestCircular=<yes|no>
```

\*covStat: Is the log of the ratio of the contig being unique vs it being two-copy, based on the number of reads. Positive means more likely to be unique, negative means more likely to be repetitive

# Canu output files - all found in helpdocs

## GRAPHS

- **ecoli\*.gkpStore** - \*.svg

## LAYOUT

The layout provides information on where each read ended up in the final assembly, including contig and positions. It also includes the consensus sequence for each contig.

- **ecoli\*.unitigs.layout.readToTig** - The position of each read in a contig
- **ecoli\*.contigs.layout.tigInfo**, **ecoli\*.unitigs.layout.tigInfo** - A list of the contigs (unitigs), lengths, coverage, number of reads and other metadata.

# Genome coverage

- How much coverage is enough?  
4.5 Mb - 223Mb fastq reads (25X), 13,124 reads

Try some assemblies with filtered data sets.



# Using canu - split reads

- Count and split files:

```
$ wc -l pacbio.fastq
```

```
$ split -l output number of lines - full: 52,496
```

```
- 60% 31,498
```

```
$ split -l 31498 pacbio.fasta
```

```
$ mv xaa ecoli_filtered_0.60.fastq
```

```
$ rm the rest (xa*)
```

- confirm number of lines and reads

# Genome coverage

- How much coverage is enough?  
4.5 Mb - 223Mb fastq reads (25X), 13,124 reads

Try some assemblies with filtered data sets:

0.75 - 9,843 reads, 168 Mb, ~18X

0.60 - 7,874 reads, 133Mb, ~15 X

0.50 - 6,562 reads, 111Mb, ~12.5X

0.25 - 3,281 reads, 56Mb, ~ 6.3X

# Genome coverage

- How much coverage is enough?

Full data set:

4.8 Mb - 223Mb fastq reads (25X), 13,124 reads

Percent Reads	Coverage	Contigs	Unitigs	Unassembled Reads	Genome Size (Mb)
1	25X	1	1	1,506	4.6
0.75	18X	1	1	1,109	4.6
0.6	15X	8	8	959	4.6
0.5	12.5X	33	35	996	4.5
0.25	6.3X	67	75	1,188	1.9

[https://rtsf.natsci.msu.edu/\\_rtsf/assets/File/depth%20and%20coverage.pdf](https://rtsf.natsci.msu.edu/_rtsf/assets/File/depth%20and%20coverage.pdf)

# Genome coverage

## SEQUENCE

- **ecoli\*.contigs.fasta** - Everything which could be assembled and is part of the primary assembly, including both unique and repetitive elements.
- **ecoli\*.unitigs.fasta** - Contigs, split at alternate paths in the graph
- **ecoli\*.unassembled.fasta** - Reads and low-coverage contigs which could not be incorporated into the primary assembly.

The header line for each sequence provides some metadata on the sequence:

```
tig # len=<integer> reads=<integer> covStat=<float> gappedBases=<yes|no>  
class=<contig|bubble|unassm> suggestRepeat=<yes|no> suggestCircular=<yes|no>
```

**ecoli-0.25.contigs.fasta:**

```
>tig00000001 len=102939 reads=69 covStat=1.36 gappedBases=no class=contig  
suggestRepeat=no suggestCircular=no
```

**ecoli-full.contigs.fasta:**

```
>tig00000001 len=4639282 reads=10263 covStat=3780.63 gappedBases=no  
class=contig suggestRepeat=no suggestCircular=no
```

# Poisson Distribution

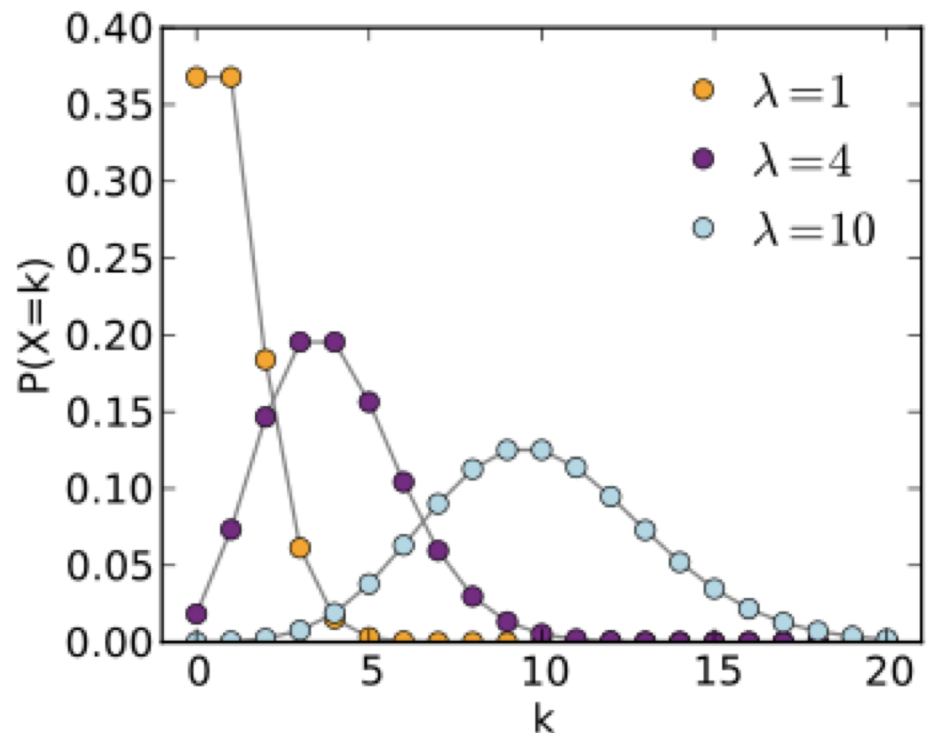
The probability of a given number of events occurring in a fixed interval of time and/or space if these events occur with a known average rate and independently of the time since the last event.

Resembles a normal distribution, but over the positive values, and with only a single parameter.

**Key property:**

- ***The standard deviation is the square root of the mean.***

$$P(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$



# Genome Assembly Problem Set

Using `ecoli-0.25.contigs.fasta`, write a script that reports:

1. The number of contigs in the file
2. The shortest contig.
3. The longest contig.
4. The L50 size
5. The N50 size

The `fasta` file is listed on [github workshops/GenomeAssembly:ecoliPB-filtered\\_0.50.contigs.fasta](#).

# N50 size

50% of the genome in contigs as long as or larger than N50 value

Example: 1 Mbp genome

50%



N50 size = 30 kbp

$(300k + 100k + 45k + 45k + 30k = 520k \geq 500kbp)$

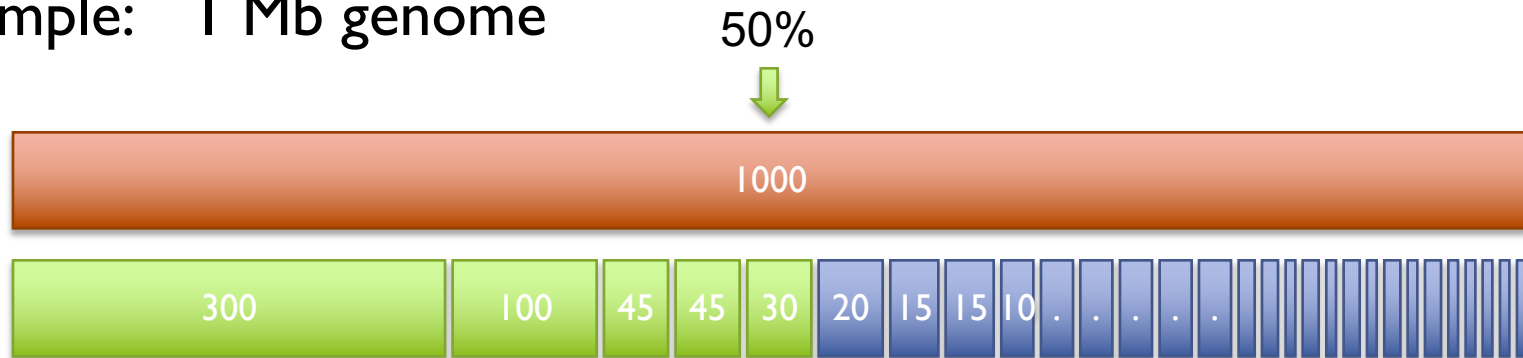
**A greater N50 is usually a sign of assembly improvement**

- Comparable with genomes of similar size
- Genome composition can bias comparisons
- High L50 vs Low N50

# L50 size

Number of contigs that are as long or longer than the N50 value

Example: 1 Mb genome



L50 size = 30 kbp

$(300k + 100k + 45k + 45k + 30k = 520k \geq 500kbp)$

- High L50 vs Low N50
  - longer sequences and fewer of them....in theory
  - lower stringency can inflate N50