

---

# **canu Documentation**

***Release 1.8***

**Adam Phillippy, Sergey Koren, Brian Walenz**

**Oct 22, 2018**



<b>1</b>	<b>Canu Quick Start</b>	<b>1</b>
1.1	Assembling PacBio or Nanopore data . . . . .	1
1.2	Assembling With Multiple Technologies and Multiple Files . . . . .	2
1.3	Assembling Low Coverage Datasets . . . . .	3
1.4	Trio Binning Assembly . . . . .	3
1.5	Consensus Accuracy . . . . .	4
<b>2</b>	<b>Canu FAQ</b>	<b>5</b>
2.1	What resources does Canu require for a bacterial genome assembly? A mammalian assembly? . . . . .	6
2.2	How do I run Canu on my SLURM / SGE / PBS / LSF / Torque system? . . . . .	6
2.3	My run stopped with the error 'Failed to submit batch jobs' . . . . .	6
2.4	My run of Canu was killed by the sysadmin; the power going out; my cat stepping on the power button; et cetera. Is it safe to restart? How do I restart? . . . . .	7
2.5	My genome size and assembly size are different, help! . . . . .	7
2.6	What parameters should I use for my reads? . . . . .	7
2.7	Can I assemble RNA sequence data? . . . . .	8
2.8	My assembly is running out of space, is too slow? . . . . .	8
2.9	My assembly continuity is not good, how can I improve it? . . . . .	9
2.10	What parameters can I tweak? . . . . .	9
2.11	My asm.contigs.fasta is empty, why? . . . . .	11
2.12	Why is my assembly missing my favorite short plasmid? . . . . .	11
2.13	Why do I get less corrected read data than I asked for? . . . . .	12
2.14	What is the minimum coverage required to run Canu? . . . . .	12
2.15	My circular element is duplicated/has overlap? . . . . .	12
2.16	My genome is AT (or GC) rich, do I need to adjust parameters? What about highly repetitive genomes? . . . . .	12
2.17	How can I send data to you? . . . . .	12
<b>3</b>	<b>Canu Tutorial</b>	<b>15</b>
3.1	Canu, the command . . . . .	15
3.2	Canu, the pipeline . . . . .	16
3.3	Module Tags . . . . .	16
3.4	Execution Configuration . . . . .	17
3.5	Error Rates . . . . .	18
3.6	Minimum Lengths . . . . .	19
3.7	Overlap configuration . . . . .	19
3.8	Ovl Overlapper Configuration . . . . .	19
3.9	Ovl Overlapper Parameters . . . . .	19

3.10	Mhap Overlapper Parameters . . . . .	20
3.11	Minimap Overlapper Parameters . . . . .	20
3.12	Outputs . . . . .	21
<b>4</b>	<b>Canu Pipeline</b>	<b>25</b>
<b>5</b>	<b>Canu Parameter Reference</b>	<b>27</b>
5.1	Global Options . . . . .	27
5.2	Process Control . . . . .	28
5.3	General Options . . . . .	29
5.4	File Staging . . . . .	30
5.5	Cleanup Options . . . . .	30
5.6	Executive Configuration . . . . .	31
5.7	Overlapper Configuration . . . . .	31
5.8	Overlap Store . . . . .	32
5.9	Meryl . . . . .	33
5.10	Overlap Based Trimming . . . . .	33
5.11	Grid Engine Support . . . . .	33
5.12	Algorithm Selection . . . . .	35
5.13	Overlap Error Adjustment . . . . .	36
5.14	Unitigger . . . . .	37
5.15	Consensus Partitioning . . . . .	37
5.16	Read Correction . . . . .	37
5.17	Output Filtering . . . . .	38
<b>6</b>	<b>Canu Command Reference</b>	<b>39</b>
<b>7</b>	<b>Software Background</b>	<b>41</b>
7.1	References . . . . .	41
<b>8</b>	<b>Publication</b>	<b>43</b>
<b>9</b>	<b>Install</b>	<b>45</b>
<b>10</b>	<b>Learn</b>	<b>47</b>

---

## Canu Quick Start

---

Canu specializes in assembling PacBio or Oxford Nanopore sequences. Canu operates in three phases: correction, trimming and assembly. The correction phase will improve the accuracy of bases in reads. The trimming phase will trim reads to the portion that appears to be high-quality sequence, removing suspicious regions such as remaining SMRTbell adapter. The assembly phase will order the reads into contigs, generate consensus sequences and create graphs of alternate paths.

For eukaryotic genomes, coverage more than 20x is enough to outperform current hybrid methods, however, between 30x and 60x coverage is the recommended minimum. More coverage will let Canu use longer reads for assembly, which will result in better assemblies.

Input sequences can be FASTA or FASTQ format, uncompressed or compressed with gzip (.gz), bzip2 (.bz2) or xz (.xz). Note that zip files (.zip) are not supported.

Canu can resume incomplete assemblies, allowing for recovery from system outages or other abnormal terminations. On each restart of Canu, it will examine the files in the assembly directory to decide what to do next. For example, if all but two overlap tasks have finished, only the two that are missing will be computed. For best results, do not change Canu parameters between restarts.

Canu will auto-detect computational resources and scale itself to fit, using all of the resources available and are reasonable for the size of your assembly. Memory and processors can be explicitly limited with with parameters *maxMemory* and *maxThreads*. See section *Execution Configuration* for more details.

Canu will automatically take full advantage of any LSF/PBS/PBSPPro/Torque/Slurm/SGE grid available, even submitting itself for execution. Canu makes heavy use of array jobs and requires job submission from compute nodes, which are sometimes not available or allowed. Canu option `useGrid=false` will restrict Canu to using only the current machine, while option `useGrid=remote` will configure Canu for grid execution but not submit jobs to the grid. See section *Execution Configuration* for more details.

The *Canu Tutorial* has more background, and the *Canu FAQ* has a wealth of practical advice.

## 1.1 Assembling PacBio or Nanopore data

Pacific Biosciences released P6-C4 chemistry reads for Escherichia coli K12. You can [download them from their original release](#), but note that you must have the [SMRTpipe software](#) installed to extract the reads as FASTQ. Instead,

use a [FASTQ format 25X subset](#) (223MB). Download from the command line with:

```
curl -L -o pacbio.fastq http://gembox.cbcb.umd.edu/mhap/raw/ecoli_p6_25x.filtered.  
↪fastq
```

There doesn't appear to be any "official" Oxford Nanopore sample data, but the [Loman Lab](#) released a set of runs, also for Escherichia coli K12. This is early data, from September 2015. Any of the four runs will work; we picked [MAP-006-1](#) (243 MB). Download from the command line with:

```
curl -L -o oxford.fasta http://nanopore.s3.climb.ac.uk/MAP006-PCR-1_2D_pass.fasta
```

By default, Canu will correct the reads, then trim the reads, then assemble the reads to unitigs. Canu needs to know the approximate genome size (so it can determine coverage in the input reads) and the technology used to generate the reads.

For PacBio:

```
canu \  
-p ecoli -d ecoli-pacbio \  
genomeSize=4.8m \  
-pacbio-raw pacbio.fastq
```

For Nanopore:

```
canu \  
-p ecoli -d ecoli-oxford \  
genomeSize=4.8m \  
-nanopore-raw oxford.fasta
```

Output and intermediate files will be in directories 'ecoli-pacbio' and 'ecoli-nanopore', respectively. Intermediate files are written in directories 'correction', 'trimming' and 'unitigging' for the respective stages. Output files are named using the '-p' prefix, such as 'ecoli.contigs.fasta', 'ecoli.contigs.gfa', etc. See section [Outputs](#) for more details on outputs (intermediate files aren't documented).

## 1.2 Assembling With Multiple Technologies and Multiple Files

Canu can use reads from any number of input files, which can be a mix of formats and technologies. We'll assemble a mix of 10X PacBio reads in two FASTQ files and 10X of Nanopore reads in one FASTA file:

```
curl -L -o mix.tar.gz http://gembox.cbcb.umd.edu/mhap/raw/ecoliP6Oxford.tar.gz  
tar xvzf mix.tar.gz  
  
canu \  
-p ecoli -d ecoli-mix \  
genomeSize=4.8m \  
-pacbio-raw pacbio.part?.fastq.gz \  
-nanopore-raw oxford.fasta.gz
```

### 1.2.1 Correct, Trim and Assemble, Manually

Sometimes, however, it makes sense to do the three top-level tasks by hand. This would allow trying multiple unitig construction parameters on the same set of corrected and trimmed reads, or skipping trimming and assembly if you only want corrected reads.

We'll use the PacBio reads from above. First, correct the raw reads:

```
canu -correct \
  -p ecoli -d ecoli \
  genomeSize=4.8m \
  -pacbio-raw pacbio.fastq
```

Then, trim the output of the correction:

```
canu -trim \
  -p ecoli -d ecoli \
  genomeSize=4.8m \
  -pacbio-corrected ecoli/ecoli.correctedReads.fasta.gz
```

And finally, assemble the output of trimming, twice, with different stringency on which overlaps to use (see *correctedErrorRate*):

```
canu -assemble \
  -p ecoli -d ecoli-erate-0.039 \
  genomeSize=4.8m \
  correctedErrorRate=0.039 \
  -pacbio-corrected ecoli/ecoli.trimmedReads.fasta.gz

canu -assemble \
  -p ecoli -d ecoli-erate-0.075 \
  genomeSize=4.8m \
  correctedErrorRate=0.075 \
  -pacbio-corrected ecoli/ecoli.trimmedReads.fasta.gz
```

Note that the assembly stages use different ‘-d’ directories. It is not possible to run multiple copies of canu with the same work directory.

## 1.3 Assembling Low Coverage Datasets

We claimed Canu works down to 20X coverage, and we will now assemble a 20X subset of *S. cerevisiae* (215 MB). When assembling, we adjust *correctedErrorRate* to accommodate the slightly lower quality corrected reads:

```
curl -L -o yeast.20x.fastq.gz http://gembox.cbcb.umd.edu/mhap/raw/yeast_filtered.20x.
↪fastq.gz

canu \
  -p asm -d yeast \
  genomeSize=12.1m \
  correctedErrorRate=0.105 \
  -pacbio-raw yeast.20x.fastq.gz
```

## 1.4 Trio Binning Assembly

Canu has support for using parental short-read sequencing to classify and bin the F1 reads (see [Trio Binning manuscript](#) for details). This example demonstrates the functionality using a synthetic mix of two *Escherichia coli* datasets. First download the data:

```
curl -L -o K12.parental.fasta https://gembox.cbcb.umd.edu/triobinning/example/k12.12.
↪fasta
curl -L -o O157.parental.fasta https://gembox.cbcb.umd.edu/triobinning/example/o157.
↪12.fasta
curl -L -o F1.fasta https://gembox.cbcb.umd.edu/triobinning/example/pacbio.fasta

trioCanu \
-p asm -d ecoliTrio \
genomeSize=5m \
-haplotypeK12 K12.parental.fasta \
-haplotypeO157 O157.parental.fasta \
-pacbio-raw F1.fasta
```

The run will produce two assemblies, `ecoliTrio/haplotypeK12/asm.contigs.fasta` and `ecoliTrio/haplotypeO157/asm.contigs.fasta`. As comparison, you can try co-assembling the datasets instead:

```
canu \
-p asm -d ecoliHap \
genomeSize=5m \
corOutCoverage=200 "batOptions=-dg 3 -db 3 -dr 1 -ca 500 -cp 50" \
-pacbio-raw F1.fasta
```

and compare the contiguity/accuracy. The current version of `trioCanu` is not yet optimized for memory use so requires adjusted parameters for large genomes. Adding the options:

```
gridOptionsExecutive="--mem=250g" gridOptionsMeryl='--partition=largemem --mem=1000g'
```

should be sufficient for a mammalian genome.

## 1.5 Consensus Accuracy

Canu consensus sequences are typically well above 99% identity for PacBio datasets. Nanopore accuracy varies depending on pore and basecaller version, but is typically above 98% for recent data. Accuracy can be improved by polishing the contigs with tools developed specifically for that task. We recommend [Quiver](#) for PacBio and [Nanopolish](#) for Oxford Nanopore data. When Illumina reads are available, [Pilon](#) can be used to polish either PacBio or Oxford Nanopore assemblies.



- *What resources does Canu require for a bacterial genome assembly? A mammalian assembly?*
- *How do I run Canu on my SLURM / SGE / PBS / LSF / Torque system?*
- *My run stopped with the error 'Failed to submit batch jobs'*
- *My run of Canu was killed by the sysadmin; the power going out; my cat stepping on the power button; et cetera. Is it safe to restart? How do I restart?*
- *My genome size and assembly size are different, help!*
- *What parameters should I use for my reads?*
- *Can I assemble RNA sequence data?*
- *My assembly is running out of space, is too slow?*
- *My assembly continuity is not good, how can I improve it?*
- *What parameters can I tweak?*
- *My asm.contigs.fasta is empty, why?*
- *Why is my assembly missing my favorite short plasmid?*
- *Why do I get less corrected read data than I asked for?*
- *What is the minimum coverage required to run Canu?*
- *My circular element is duplicated/has overlap?*
- *My genome is AT (or GC) rich, do I need to adjust parameters? What about highly repetitive genomes?*
- *How can I send data to you?*

## 2.1 What resources does Canu require for a bacterial genome assembly? A mammalian assembly?

Canu will detect available resources and configure itself to run efficiently using those resources. It will request resources, for example, the number of compute threads to use, Based on the genome size being assembled. It will fail to even start if it feels there are insufficient resources available.

A typical bacterial genome can be assembled with 8GB memory in a few CPU hours - around an hour on 8 cores. It is possible, but not allowed by default, to run with only 4GB memory.

A well-behaved large genome, such as human or other mammals, can be assembled in 10,000 to 25,000 CPU hours, depending on coverage. A grid environment is strongly recommended, with at least 16GB available on each compute node, and one node with at least 64GB memory. You should plan on having 3TB free disk space, much more for highly repetitive genomes.

Our compute nodes have 48 compute threads and 128GB memory, with a few larger nodes with up to 1TB memory. We develop and test (mostly bacteria, yeast and drosophila) on laptops and desktops with 4 to 12 compute threads and 16GB to 64GB memory.

## 2.2 How do I run Canu on my SLURM / SGE / PBS / LSF / Torque system?

Canu will detect and configure itself to use on most grids. Canu will NOT request explicit time limits or queues/partitions. You can supply your own grid options, such as a partition on SLURM, an account code on SGE, and/or time limits with `gridOptions=<your options list>` which will be passed to every job submitted by Canu. Similar options exist for every stage of Canu, which could be used to, for example, restrict overlapping to a specific partition or queue.

To disable grid support and run only on the local machine, specify `useGrid=false`

It is possible to limit the number of grid jobs running at the same time, but this isn't directly supported by Canu. The various `gridOptions` parameters can pass grid-specific parameters to the submit commands used; see [Issue #756](#) for Slurm and SGE examples.

## 2.3 My run stopped with the error 'Failed to submit batch jobs'

The grid you run on must allow compute nodes to submit jobs. This means that if you are on a compute host, `qsub`/`bsub`/`sbatch`/`etc` must be available and working. You can test this by starting an interactive compute session and running the submit command manually (e.g. `qsub` on SGE, `bsub` on LSF, `sbatch` on SLURM).

If this is not the case, Canu **WILL NOT** work on your grid. You must then set `useGrid=false` and run on a single machine. Alternatively, you can run Canu with `useGrid=remote` which will stop at every submit command, list what should be submitted. You then submit these jobs manually, wait for them to complete, and run the Canu command again. This is a manual process but currently the only workaround for grids without submit support on the compute nodes.

## 2.4 My run of Canu was killed by the sysadmin; the power going out; my cat stepping on the power button; et cetera. Is it safe to restart? How do I restart?

Yes, perfectly safe! It's actually how Canu runs normally: each time Canu starts, it examines the state of the assembly to decide what it should do next. For example, if six overlap tasks have no results, it'll run just those six tasks.

This also means that if you want to redo some step, just remove those results from the assembly directory. Some care needs to be taken to make sure results computed after those are also removed.

Short answer: just rerun the `_exact_` same command as before. It'll do the right thing.

## 2.5 My genome size and assembly size are different, help!

The difference could be due to a heterozygous genome where the assembly separated some loci. It could also be because the previous estimate is incorrect. We typically use two analyses to see what happened. First, a BUSCO analysis will indicate duplicated genes. For example this assembly:

```
INFO      C:98.5%[S:97.9%,D:0.6%],F:1.0%,M:0.5%,n:2799
INFO      2756 Complete BUSCOs (C)
INFO      2740 Complete and single-copy BUSCOs (S)
INFO      16 Complete and duplicated BUSCOs (D)
```

does not have much duplication but this assembly:

```
INFO      C:97.6%[S:15.8%,D:81.8%],F:0.9%,M:1.5%,n:2799
INFO      2732 Complete BUSCOs (C)
INFO      443 Complete and single-copy BUSCOs (S)
INFO      2289 Complete and duplicated BUSCOs (D)
```

does. We have had some success (in limited testing) using `purge_haplotigs` to remove duplication. `Purge haplotigs` will also generate a coverage plot which will usually have two peaks when assemblies have separated some loci.

## 2.6 What parameters should I use for my reads?

Canu is designed to be universal on a large range of PacBio (C2, P4-C2, P5-C3, P6-C4) and Oxford Nanopore (R6 through R9) data. Assembly quality and/or efficiency can be enhanced for specific datatypes:

**Nanopore R7 1D and Low Identity Reads** With R7 1D sequencing data, and generally for any raw reads lower than 80% identity, five to ten rounds of error correction are helpful:

```
canu -p r1 -d r1 -correct corOutCoverage=500 corMinCoverage=0_
↳corMhapSensitivity=high -nanopore-raw your_reads.fasta
canu -p r2 -d r2 -correct corOutCoverage=500 corMinCoverage=0_
↳corMhapSensitivity=high -nanopore-raw r1/r1.correctedReads.
↳fasta.gz
canu -p r3 -d r3 -correct corOutCoverage=500 corMinCoverage=0_
↳corMhapSensitivity=high -nanopore-raw r2/r2.correctedReads.
↳fasta.gz
```

(continues on next page)

(continued from previous page)

```

canu -p r4 -d r4 -correct corOutCoverage=500 corMinCoverage=0
↳corMhapSensitivity=high -nanopore-raw r3/r3.correctedReads.
↳fasta.gz
canu -p r5 -d r5 -correct corOutCoverage=500 corMinCoverage=0
↳corMhapSensitivity=high -nanopore-raw r4/r4.correctedReads.
↳fasta.gz

```

Then assemble the output of the last round, allowing up to 30% difference in overlaps:

```

canu -p asm -d asm correctedErrorRate=0.3 utgGraphDeviation=50 -
↳nanopore-corrected r5/r5.correctedReads.fasta.gz

```

**Nanopore R7 2D and Nanopore R9 1D** The defaults were designed with these datasets in mind so they should work. Having very high coverage or very long Nanopore reads can slow down the assembly significantly. You can try the `overlapper=mhap utgReAlign=true` option which is much faster but may produce less contiguous assemblies on large genomes.

**Nanopore R9 2D and PacBio P6** Slightly decrease the maximum allowed difference in overlaps from the default of 14.4% to 12.0% with `correctedErrorRate=0.120`

#### PacBio Sequel

Based on an *A. thaliana* dataset, and a few more recent mammalian genomes, slightly increase the maximum allowed difference from the default of 4.5% to 8.5% with `correctedErrorRate=0.085 corMhapSensitivity=normal`.

Only add the second parameter (`corMhapSensitivity=normal`) if you have >50x coverage.

**Nanopore R9 large genomes** Due to some systematic errors, the identity estimate used by Canu for correction can be an over-estimate of true error, inflating runtime. For recent large genomes (>1gpb) with more than 30x coverage, we've used '`corMhapOptions=--threshold 0.8 --num-hashes 512 --ordered-sketch-size 1000 --ordered-kmer-size 14`'. This is not needed for below 30x coverage.

## 2.7 Can I assemble RNA sequence data?

Canu will likely mis-assemble, or completely fail to assemble, RNA data. It will do a reasonable job at generating corrected reads though. Reads are corrected using (local) best alignments to other reads, and alignments between different isoforms are usually obviously not 'best'. Just like with DNA sequences, similar isoforms can get 'mixed' together. We've heard of reasonable success from users, but do not have any parameter suggestions to make.

Note that Canu will silently translate 'U' bases to 'T' bases on input, but **NOT** translate the output bases back to 'U'.

## 2.8 My assembly is running out of space, is too slow?

We don't have a good way to estimate of disk space used for the assembly. It varies with genome size, repeat content, and sequencing depth. A human genome sequenced with PacBio or Nanopore at 40-50x

typically requires 1-2TB of space at the peak. Plants, unfortunately, seem to want a lot of space. 10TB is a reasonable guess. We've seen it as bad as 20TB on some very repetitive genomes.

The most common cause of high disk usage is a very repetitive or large genome. There are some parameters you can tweak to both reduce disk space and speed up the run. Try adding the options `corMhapFilterThreshold=0.0000000002 corMhapOptions="--threshold 0.80 --num-hashes 512 --num-min-matches 3 --ordered-sketch-size 1000 --ordered-kmer-size 14 --min-olap-length 2000 --repeat-idf-scale 50" mhapMemory=60g mhapBlockSize=500 ovlMerThreshold=500`. This will suppress repeats more than the default settings and speed up both correction and assembly.

It is also possible to clean up some intermediate outputs before the assembly is complete to save space. If you already have a `*.ovlStore.BUILDING/1-bucketize.successs` file in your current step (e.g. `correct``), you can clean up the files under `1-overlapper/blocks`. You can also remove the `ovlStore` for the previous step if you have its output (e.g. if you have `asm.trimmedReads.fasta`, `gz`, you can remove `trimming/asm.ovlStore`).

## 2.9 My assembly continuity is not good, how can I improve it?

The most important determinant for assembly quality is sequence length, followed by the repeat complexity/heterozygosity of your sample. The first thing to check is the amount of corrected bases output by the correction step. This is logged in the stdout of Canu or in `canu-scripts/canu.*.out` if you are running in a grid environment. For example on a haploid *H. sapiens* sample:

```

-- BEGIN TRIMMING
--
...
-- In gatekeeper store 'chm1/trimming/asm.gkpStore':
--   Found 5459105 reads.
--   Found 91697412754 bases (29.57 times coverage).
...

```

Canu tries to correct the longest 40X of data. Some loss is normal but having output coverage below 20-25X is a sign that correction did not work well (assuming you have more input coverage than that). If that is the case, re-running with `corMhapSensitivity=normal` if you have >50X or `corMhapSensitivity=high corMinCoverage=0` otherwise can help. You can also increase the target coverage to correct `corOutCoverage=100` to get more correct sequences for assembly. If there are sufficient corrected reads, the poor assembly is likely due to either repeats in the genome being greater than read lengths or a high heterozygosity in the sample. Stay tuned for mor information on tuning unitigging in those instances.

## 2.10 What parameters can I tweak?

For all stages:

- `rawErrorRate` is the maximum expected difference in an alignment of two `_uncorrected_` reads. It is a meta-parameter that sets other parameters.
- `correctedErrorRate` is the maximum expected difference in an alignment of two `_corrected_` reads. It is a meta-parameter that sets other parameters. (If you're used to the `errorRate` parameter, multiply that by 3 and use it here.)

- `minReadLength` and `minOverlapLength`. The defaults are to discard reads shorter than 1000bp and to not look for overlaps shorter than 500bp. Increasing `minReadLength` can improve run time, and increasing `minOverlapLength` can improve assembly quality by removing false overlaps. However, increasing either too much will quickly degrade assemblies by either omitting valuable reads or missing true overlaps.

For correction:

- `corOutCoverage` controls how much coverage in corrected reads is generated. The default is to target 40X, but, for various reasons, this results in 30X to 35X of reads being generated.
- `corMinCoverage`, loosely, controls the quality of the corrected reads. It is the coverage in evidence reads that is needed before a (portion of a) corrected read is reported. Corrected reads are generated as a consensus of other reads; this is just the minimum coverage needed for the consensus sequence to be reported. The default is based on input read coverage: 0x coverage for less than 30X input coverage, and 4x coverage for more than that.

For assembly:

- `utgOvlErrorRate` is essentially a speed optimization. Overlaps above this error rate are not computed. Setting it too high generally just wastes compute time, while setting it too low will degrade assemblies by missing true overlaps between lower quality reads.
- `utgGraphDeviation` and `utgRepeatDeviation` what quality of overlaps are used in contig construction or in breaking contigs at false repeat joins, respectively. Both are in terms of a deviation from the mean error rate in the longest overlaps.
- `utgRepeatConfusedBP` controls how similar a true overlap (between two reads in the same contig) and a false overlap (between two reads in different contigs) need to be before the contig is split. When this occurs, it isn't clear which overlap is 'true' - the longer one or the slightly shorter one - and the contig is split to avoid misassemblies.

For polyploid genomes:

Generally, there's a couple of ways of dealing with the ploidy.

1. **Avoid collapsing the genome** so you end up with double (assuming diploid) the genome size as long as your divergence is above about 2% (for PacBio data). Below this divergence, you'd end up collapsing the variations. We've used the following parameters for polyploid populations (PacBio data):

```
corOutCoverage=200 "batOptions=-dg 3 -db 3 -dr 1 -ca 500
-cp 50"
```

This will output more corrected reads (than the default 40x). The latter option will be more conservative at picking the error rate to use for the assembly to try to maintain haplotype separation. If it works, you'll end up with an assembly  $\geq 2x$  your haploid genome size. Post-processing using gene information or other synteny information is required to remove redundancy from this assembly.

2. **Smash haplotypes together** and then do phasing using another approach (like HapCUT2 or whatshap or others). In that case you want to do the opposite, increase the error rates used for finding overlaps:

```
corOutCoverage=200 correctedErrorRate=0.15
```

When trimming, reads will be trimmed using other reads in the same chromosome (and probably some reads from other chromosomes). When assembling, overlaps well outside the observed error rate distribution are discarded.

We typically prefer option 1 which will lead to a larger than expected genome size. We have had some success (in limited testing) using `purge_haplotigs` to remove this duplication.

For metagenomes:

The basic idea is to use all data for assembly rather than just the longest as default. The parameters we've used recently are:

```
corOutCoverage=10000 corMhapSensitivity=high
corMinCoverage=0 redMemory=32 oeaMemory=32 batMemory=200
```

For low coverage:

- For less than 30X coverage, increase the allowed difference in overlaps by a few percent (from 4.5% to 8.5% (or more) with `correctedErrorRate=0.105` for PacBio and from 14.4% to 16% (or more) with `correctedErrorRate=0.16` for Nanopore), to adjust for inferior read correction. Canu will automatically reduce `corMinCoverage` to zero to correct as many reads as possible.

For high coverage:

- For more than 60X coverage, decrease the allowed difference in overlaps (from 4.5% to 4.0% with `correctedErrorRate=0.040` for PacBio, from 14.4% to 12% with `correctedErrorRate=0.12` for Nanopore), so that only the better corrected reads are used. This is primarily an optimization for speed and generally does not change assembly continuity.

## 2.11 My asm.contigs.fasta is empty, why?

Canu creates three assembled sequence *output files*: `<prefix>.contigs.fasta`, `<prefix>.unitigs.fasta`, and `<prefix>.unassembled.fasta`, where contigs are the primary output, unitigs are the primary output split at alternate paths, and unassembled are the leftover pieces.

The `contigFilter` parameter sets several parameters that control how small or low coverage initial contigs are handled. By default, initial contigs with more than 50% of the length at less than 3X coverage will be classified as 'unassembled' and removed from the assembly, that is, `contigFilter="2 0 1.0 0.5 3"`. The filtering can be disabled by changing the last number from '3' to '0' (meaning, filter if 50% of the contig is less than 0X coverage).

## 2.12 Why is my assembly is missing my favorite short plasmid?

In Canu v1.6 and earlier only the longest 40X of data (based on the specified genome size) is used for correction. Datasets with uneven coverage or small plasmids can fail to generate enough corrected reads to give enough coverage for assembly, resulting in gaps in the genome or even no reads for small plasmids. Set `corOutCoverage=1000` (or any value greater than your total input coverage) to correct all input data.

An alternate approach is to correct all reads (`-correct corOutCoverage=1000`) then assemble 40X of reads picked at random from the `<prefix>.correctedReads.fasta.gz` output.

More recent Canu versions dynamically select poorly represented sequences to avoid missing short plasmids so this should no longer happen.

## 2.13 Why do I get less corrected read data than I asked for?

Some reads are trimmed during correction due to being chimeric or because there wasn't enough evidence to generate a quality corrected sequence. Typically, this results in a 25% loss. Setting `corMinCoverage=0` will report all bases, even low those of low quality. Canu will trim these in its 'trimming' phase before assembly.

## 2.14 What is the minimum coverage required to run Canu?

For eukaryotic genomes, coverage more than 20X is enough to outperform current hybrid methods. Below that, you will likely not assemble the full genome.

## 2.15 My circular element is duplicated/has overlap?

This is expected for any circular elements. They can overlap by up to a read length due to how Canu constructs contigs. Canu provides an alignment string in the GFA output which can be converted to an alignment to identify the trimming points.

An alternative is to run MUMmer to get self-alignments on the contig and use those trim points. For example, assuming the circular element is in `tig00000099.fa`. Run:

```
nucmer -maxmatch -nosimplify tig00000099.fa tig00000099.fa
show-coords -lrcTH out.delta
```

to find the end overlaps in the tig. The output would be something like:

```
1 1895    48502   50400   1895    1899    99.37   50400   50400   3.76    3.
↪77    tig00000001    tig00000001
48502   50400   1        1895    1899    99.37   50400   50400   3.
↪77    3.76    tig00000001    tig00000001
```

means trim to 1 to 48502. There is also an alternate `writeup`.

## 2.16 My genome is AT (or GC) rich, do I need to adjust parameters? What about highly repetitive genomes?

On bacterial genomes, no adjustment of parameters is (usually) needed. See the next question.

On repetitive genomes with with a significantly skewed AT/GC ratio, the Jaccard estimate used by MHAP is biased. Setting `corMaxEvidenceErate=0.15` is sufficient to correct for the bias in our testing.

In general, with high coverage repetitive genomes (such as plants) it can be beneficial to set the above parameter anyway, as it will eliminate repetitive matches, speed up the assembly, and sometime improve unitigs.

## 2.17 How can I send data to you?

FTP to <ftp://ftp.cbcb.umd.edu/incoming/sergek>. This is a write-only location that only the Canu developers can see.



Here is a quick walk-through using a command-line ftp client (should be available on most Linux and OSX installations). Say we want to transfer a file named `reads.fastq`. First, run `ftp ftp.cbc.umd.edu`, specify `anonymous` as the user name and hit return for password (blank). Then `cd incoming/sergek`, `put reads.fastq`, and `quit`.

That's it, you won't be able to see the file but we can download it.



Canu assembles reads from PacBio RS II or Oxford Nanopore MinION instruments into uniquely-assemblable contigs, unitigs. Canu owes lots of its design and code to [celera-assembler](#).

Canu can be run using hardware of nearly any shape or size, anywhere from laptops to computational grids with thousands of nodes. Obviously, larger assemblies will take a long time to compute on laptops, and smaller assemblies can't take advantage of hundreds of nodes, so what is being assembled plays some part in determining what hardware can be effectively used.

Most algorithms in canu have been multi-threaded (to use all the cores on a single node), parallelized (to use all the nodes in a grid), or both (all the cores on all the nodes).

### 3.1 Canu, the command

The **canu** command is the 'executive' program that runs all modules of the assembler. It oversees each of the three top-level tasks (correction, trimming, unitig construction), each of which consists of many steps. Canu ensures that input files for each step exist, that each step successfully finished, and that the output for each step exists. It does minor bits of processing, such as reformatting files, but generally just executes other programs.

```
canu [-correct | -trim | -assemble | -trim-assemble] \  
  [-s <assembly-specifications-file>] \  
  -p <assembly-prefix> \  
  -d <assembly-directory> \  
  genomeSize=<number>[g|m|k] \  
  [other-options] \  
  [-pacbio-raw | -pacbio-corrected | -nanopore-raw | -nanopore-corrected] *fastq
```

The **-p** option, to set the file name prefix of intermediate and output files, is mandatory. If **-d** is not supplied, canu will run in the current directory, otherwise, Canu will create the *assembly-directory* and run in that directory. It is not possible to run two different assemblies in the same directory.

The **-s** option will import a list of parameters from the supplied specification ('spec') file. These parameters will be applied before any from the command line are used, providing a method for setting commonly used parameters, but overriding them for specific assemblies.

By default, all three top-level tasks are performed. It is possible to run exactly one task by using the `-correct`, `-trim` or `-assemble` options. These options can be useful if you want to correct reads once and try many different assemblies. We do exactly that in the *Canu Quick Start*. Additionally, supplying pre-corrected reads with `-pacbio-corrected` or `-nanopore-corrected` will run only the trimming (`-trim`) and assembling (`-assemble`) stages.

Parameters are key=value pairs that configure the assembler. They set run time parameters (e.g., memory, threads, grid), algorithmic parameters (e.g., error rates, trimming aggressiveness), and enable or disable entire processing steps (e.g., don't correct errors, don't search for subreads). They are described later. One parameter is required: the `genomeSize` (in bases, with common SI prefixes allowed, for example, 4.7m or 2.8g; see `genomeSize`). Parameters are listed in the *Canu Parameter Reference*, but the common ones are described in this document.

Reads are supplied to canu by options that describe how the reads were generated, and what level of quality they are, for example, `-pacbio-raw` indicates the reads were generated on a PacBio RS II instrument, and have had no processing done to them. Each file of reads supplied this way becomes a 'library' of reads. The reads should have been (physically) generated all at the same time using the same steps, but perhaps sequenced in multiple batches. In canu, each library has a set of options setting various algorithmic parameters, for example, how aggressively to trim. To explicitly set library parameters, a text 'gkp' file describing the library and the input files must be created. Don't worry too much about this yet, it's an advanced feature, fully described in Section `gkp-files`.

The read-files contain sequence data in either FASTA or FASTQ format (or both! A quirk of the implementation allows files that contain both FASTA and FASTQ format reads). The files can be uncompressed, `gzip`, `bzip2` or `xz` compressed. We've found that "`gzip -1`" provides good compression that is fast to both compress and decompress. For 'archival' purposes, we use "`xz -9`".

## 3.2 Canu, the pipeline

The canu pipeline, that is, what it actually computes, comprises of computing overlaps and processing the overlaps to some result. Each of the three tasks (read correction, read trimming and unitig construction) follow the same pattern:

- Load reads into the read database, `gkpStore`.
- Compute k-mer counts in preparation for the overlap computation.
- Compute overlaps.
- Load overlaps into the overlap database, `ovlStore`.
- Do something interesting with the reads and overlaps.
  - The read correction task will replace the original noisy read sequences with consensus sequences computed from overlapping reads.
  - The read trimming task will use overlapping reads to decide what regions of each read are high-quality sequence, and what regions should be trimmed. After trimming, the single largest high-quality chunk of sequence is retained.
  - 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.

## 3.3 Module Tags

Because each of the three tasks share common algorithms (all compute overlaps, two compute consensus sequences, etc), parameters are differentiated by a short prefix 'tag' string. This lets canu have one generic parameter that can be set to different values for each stage in each task. For example, "`corOvlMemory`" will set memory usage for overlaps being generated for read correction; "`obtOvlMemory`" for overlaps generated for Overlap Based Trimming; "`utgOvlMemory`" for overlaps generated for unitig construction.

The tags are:

Tag	Usage
master	the canu script itself, and any components that it runs directly
cns	unitig consensus generation
cor	read correction generation
red	read error detection
oea	overlap error adjustment
ovl	the standard overlapper
corovl	the standard overlapper, as used in the correction phase
obtovl	the standard overlapper, as used in the trimming phase
utgovl	the standard overlapper, as used in the assembly phase
mhap	the mhap overlapper
cormhap	the mhap overlapper, as used in the correction phase
obtmhap	the mhap overlapper, as used in the trimming phase
utgmhap	the mhap overlapper, as used in the assembly phase
mmap	the <a href="#">minimap</a> overlapper
cormmap	the minimap overlapper, as used in the correction phase
obtmmap	the minimap overlapper, as used in the trimming phase
utgmmap	the minimap overlapper, as used in the assembly phase
ovb	the bucketizing phase of overlap store building
ovs	the sort phase of overlap store building

We'll get to the details eventually.

## 3.4 Execution Configuration

There are two modes that canu runs in: locally, using just one machine, or grid-enabled, using multiple hosts managed by a grid engine. LSF, PBS/Torque, PBSPro, Sun Grid Engine (and derivations), and Slurm are supported, though LSF has had limited testing. Section [Grid Engine Configuration](#) has a few hints on how to set up a new grid engine.

By default, if a grid is detected the canu pipeline will immediately submit itself to the grid and run entirely under grid control. If no grid is detected, or if option `useGrid=false` is set, canu will run on the local machine.

In both cases, Canu will auto-detect available resources and configure job sizes based on the resources and genome size you're assembling. Thus, most users should be able to run the command without modifying the defaults. Some advanced options are outlined below. Each stage has the same five configuration options, and tags are used to specialize the option to a specific stage. The options are:

**useGrid<tag>=boolean** Run this stage on the grid, usually in parallel.

**gridOptions<tag>=string** Supply this string to the grid submit command.

**<tag>Memory=integer** Use this many gigabytes of memory, per process.

**<tag>Threads** Use this many compute threads per process.

**<tag>Concurrency** If not on the grid, run this many jobs at the same time.

Global grid options, applied to every job submitted to the grid, can be set with ‘gridOptions’. This can be used to add accounting information or access credentials.

A name can be associated with this compute using ‘gridOptionsJobName’. Canu will work just fine with no name set, but if multiple canu assemblies are running at the same time, they will tend to wait for each others jobs to finish. For example, if two assemblies are running, at some point both will have overlap jobs running. Each assembly will be waiting for all jobs named ‘ovl\_asm’ to finish. Had the assemblies specified job names, gridOptionsJobName=apple and gridOptionsJobName=orange, then one would be waiting for jobs named ‘ovl\_asm\_apple’, and the other would be waiting for jobs named ‘ovl\_asm\_orange’.

### 3.5 Error Rates

Canu expects all error rates to be reported as fraction error, not as percent error. We’re not sure exactly why this is so. Previously, it used a mix of fraction error and percent error (or both!), and was a little confusing. Here’s a handy table you can print out that converts between fraction error and percent error. Not all values are shown (it’d be quite a large table) but we have every confidence you can figure out the missing values:

Fraction Error	Percent Error
0.01	1%
0.02	2%
0.03	3%
.	.
.	.
0.12	12%
.	.
.	.

Canu error rates always refer to the percent difference in an alignment of two reads, not the percent error in a single read, and not the amount of variation in your reads. These error rates are used in two different ways: they are used to limit what overlaps are generated, e.g., don’t compute overlaps that have more than 5% difference; and they are used to tell algorithms what overlaps to use, e.g., even though overlaps were computed to 5% difference, don’t trust any above 3% difference.

There are seven error rates. Three error rates control overlap creation (*corOvlErrorRate*, *obtOvlErrorRate* and *utgOvlErrorRate*), and four error rates control algorithms (*corErrorRate*, *obtErrorRate*, *utgErrorRate*, *cnsErrorRate*).

The three error rates for overlap creation apply to the *ovl* overlap algorithm and the *mhapReAlign* option used to generate alignments from *mhap* or *minimap* overlaps. Since *mhap* is used for generating correction overlaps, the *corOvlErrorRate* parameter is not used by default. Overlaps for trimming and assembling use the *ovl* algorithm, therefore, *obtOvlErrorRate* and *utgOvlErrorRate* are used.

The four algorithm error rates are used to select which overlaps can be used for correcting reads (*corErrorRate*); which overlaps can be used for trimming reads (*obtErrorRate*); which overlaps can be used for assembling reads (*utgErrorRate*). The last error rate, *cnsErrorRate*, tells the consensus algorithm to not trust read alignments above that value.

For convenience, two meta options set the error rates used with uncorrected reads (*rawErrorRate*) or used with corrected reads. (*correctedErrorRate*). The default depends on the type of read being assembled.

Parameter	PacBio	Nanopore
rawErrorRate	0.300	0.500
correctedErrorRate	0.045	0.144

In practice, only *correctedErrorRate* is usually changed. The *Canu FAQ* has *specific suggestions* on when to change this.

Canu v1.4 and earlier used the *errorRate* parameter, which set the expected rate of error in a single corrected read.

## 3.6 Minimum Lengths

Two minimum sizes are known:

**minReadLength** Discard reads shorter than this when loading into the assembler, and when trimming reads.

**minOverlapLength** Do not save overlaps shorter than this.

## 3.7 Overlap configuration

The largest compute of the assembler is also the most complicated to configure. As shown in the ‘module tags’ section, there are up to eight (!) different overlapper configurations. For each overlapper (‘ovl’ or ‘mhap’) there is a global configuration, and three specializations that apply to each stage in the pipeline (correction, trimming or assembly).

Like with ‘grid configuration’, overlap configuration uses a ‘tag’ prefix applied to each option. The tags in this instance are ‘cor’, ‘obt’ and ‘utg’.

For example:

- To change the k-mer size for all instances of the ovl overlapper, ‘merSize=23’ would be used.
- To change the k-mer size for just the ovl overlapper used during correction, ‘corMerSize=16’ would be used.
- To change the mhap k-mer size for all instances, ‘mhapMerSize=18’ would be used.
- To change the mhap k-mer size just during correction, ‘corMhapMerSize=15’ would be used.
- To use minimap for overlap computation just during correction, ‘corOverlapper=minimap’ would be used. The minimap2 executable must be symlinked from the Canu binary folder (‘Linux-amd64/bin’ or ‘Darwin-amd64/bin’ depending on your system).

## 3.8 Ovl Overlapper Configuration

**<tag>Overlapper** select the overlap algorithm to use, ‘ovl’ or ‘mhap’.

## 3.9 Ovl Overlapper Parameters

**<tag>ovlHashBlockLength** how many bases to reads to include in the hash table; directly controls process size

**<tag>ovlRefBlockSize** how many reads to compute overlaps for in one process; directly controls process time

**<tag>ovlRefBlockLength** same, but use ‘bases in reads’ instead of ‘number of reads’

**<tag>ovlHashBits** size of the hash table (SHOULD BE REMOVED AND COMPUTED, MAYBE TWO PASS)

**<tag>ovlHashLoad** how much to fill the hash table before computing overlaps (SHOULD BE REMOVED)

**<tag>ovlMerSize** size of kmer seed; smaller - more sensitive, but slower

The overlapper will not use frequent kmers to seed overlaps. These are computed by the ‘meryl’ program, and can be selected in one of three ways.

**Terminology.** A k-mer is a contiguous sequence of k bases. The read ‘ACTTA’ has two 4-mers: ACTT and CTTA. To account for reverse-complement sequence, a ‘canonical kmer’ is the lexicographically smaller of the forward and reverse-complemented kmer sequence. Kmer ACTT, with reverse complement AAGT, has a canonical kmer AAGT. Kmer CTTA, reverse-complement TAAG, has canonical kmer CTTA.

A ‘distinct’ kmer is the kmer sequence with no count associated with it. A ‘total’ kmer (for lack of a better term) is the kmer with its count. The sequence TCGTTTTTTTCGTCG has 12 ‘total’ 4-mers and 8 ‘distinct’ kmers.

TCGTTTTTTTCGTCG	count
TCGT	2 distinct-1
CGTT	1 distinct-2
GTTT	1 distinct-3
TTTT	4 distinct-4
TTTT	4 copy of distinct-4
TTTT	4 copy of distinct-4
TTTT	4 copy of distinct-4
TTTC	1 distinct-5
TTCG	1 distinct-6
TCGT	2 copy of distinct-1
CGTC	1 distinct-7
GTCG	1 distinct-8

**<tag>MerThreshold** any kmer with count higher than N is not used

**<tag>MerDistinct** pick a threshold so as to seed overlaps using this fraction of all distinct kmers in the input. In the example above, fraction 0.875 of the k-mers (7/8) will be at or below threshold 2.

**<tag>MerTotal** pick a threshold so as to seed overlaps using this fraction of all kmers in the input. In the example above, fraction 0.667 of the k-mers (8/12) will be at or below threshold 2.

**<tag>FrequentMers** don’t compute frequent kmers, use those listed in this fasta file

### 3.10 Mhap Overlapper Parameters

**<tag>MhapBlockSize** Chunk of reads that can fit into 1GB of memory. Combined with memory to compute the size of chunk the reads are split into.

**<tag>MhapMerSize** Use k-mers of this size for detecting overlaps.

**<tag>ReAlign** After computing overlaps with mhap, compute a sequence alignment for each overlap.

**<tag>MhapSensitivity** Either ‘normal’, ‘high’, or ‘fast’.

Mhap also will down-weight frequent kmers (using tf-idf), but it’s selection of frequent is not exposed.

### 3.11 Minimap Overlapper Parameters

**<tag>MMapBlockSize** Chunk of reads that can fit into 1GB of memory. Combined with memory to compute the size of chunk the reads are split into.

**<tag>MMapMerSize** Use k-mers of this size for detecting overlaps

Minimap also will ignore high-frequency minimizers, but it’s selection of frequent is not exposed.



### 3.12 Outputs

As Canu runs, it outputs status messages, execution logs, and some analysis to the console. Most of the analysis is captured in `<prefix>.report` as well.

#### LOGGING

**<prefix>.report** Most of the analysis reported during assembly. This will report the histogram of read lengths, the histogram or k-mers in the raw and corrected reads, the summary of corrected data, summary of overlaps, and the summary of contig lengths.

You can use the k-mer corrected read histograms with tools like [GenomeScope](#) to estimate heterozygosity and genome size. In particular, histograms with more than 1 peak likely indicate a heterozygous genome. See the [Canu FAQ](#) for some suggested parameters.

The corrected read report gives a summary of the fate of all input reads. The first part::

category	original raw reads w/overlaps	original raw reads w/o/overlaps
Number of Reads	250609	477
Number of Bases	2238902045	1896925
Coverage	97.344	0.082
Median	6534	2360
Mean	8933	3976
N50	11291	5756
Minimum	1012	0
Maximum	60664	41278

reports the fraction of reads which had an overlap. In this case, the majority had at least one overlap, which is good. Next:

category	evidence reads	corrected	
		raw	expected corrected
Number of Reads	229397	48006	48006
Number of Bases	2134291652	993586222	920001699
Coverage	92.795	43.199	40.000
Median	6842	15330	14106
Mean	9303	20697	19164
N50	11512	28066	26840
Minimum	1045	10184	10183
Maximum	60664	60664	59063

reports that a total of 92.8x of raw bases are candidates for correction. By default, Canu only selects the longest 40x for correction. In this case, it selects 43.2x of raw read data which it estimates will result in 40x correction. Not all raw reads survive full-length through correction:

category	rescued	
	raw	expected corrected
Number of Reads	20030	20030
Number of Bases	90137165	61903752
Coverage	3.919	2.691

(continues on next page)

(continued from previous page)

--	Median	3324	2682
--	Mean	4500	3090
--	N50	5529	3659
--	Minimum	1012	501
--	Maximum	41475	10179

The rescued reads are those which would not have contributed to the correction of the selected longest 40x subset. These could be short plasmids, mitochondria, etc. Canu includes them even though they're too short by the 40x cutoff to avoid losing sequence during assembly. Lastly:

	-----uncorrected-----	expected	
category	raw	corrected	
--	Number of Reads	183050	183050
--	Number of Bases	1157075583	951438105
--	Coverage	50.308	41.367
--	Median	5729	5086
--	Mean	6321	5197
--	N50	7467	6490
--	Minimum	0	0
--	Maximum	50522	10183

are the reads which were deemed too short to correct. If you increase `corOutCoverage`, you could get up to 41x more corrected sequence. However, unless the genome is very heterozygous, this does not typically improve the assembly and increases the running time.

The assembly statistics (NG50, etc) are reported before and after consensus calling.

## READS

**<prefix>.correctedReads.fasta.gz** The reads after correction.

**<prefix>.trimmedReads.fasta.gz** The corrected reads after overlap based trimming.

## SEQUENCE

**<prefix>.contigs.fasta** Everything which could be assembled and is the primary assembly, including both unique and repetitive elements.

**<prefix>.unitigs.fasta** Contigs, split at alternate paths in the graph.

**<prefix>.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|unasm> suggestRepeat=<yes|no> suggestCircular=<yes|no>

len
  Length of the sequence, in bp.

reads
  Number of reads used to form the contig.

covStat
  The log of the ratio of the contig being unique versus being two-copy, based on
↳the read arrival rate. Positive values indicate more likely to be unique, while
↳negative values indicate more likely to be repetitive. See `Footnote 24 <http://
↳science.sciencemag.org/content/287/5461/2196.full#ref-24>`_ in `Myers et al., A
↳Whole-Genome Assembly of Drosophila <http://science.sciencemag.org/content/287/5461/
2196.full>`_.
```

(continues on next page)

(continued from previous page)

```

gappedBases
    If yes, the sequence includes all gaps in the multialignment.

class
    Type of sequence. Unassembled sequences are primarily low-coverage sequences,
    ↳spanned by a single read.

suggestRepeat
    If yes, sequence was detected as a repeat based on graph topology or read overlaps,
    ↳to other sequences.

suggestCircular
    If yes, sequence is likely circular. Not implemented.

```

## GRAPHS

**<prefix>.contigs.gfa** Unused or ambiguous edges between contig sequences. Bubble edges cannot be represented in this format.

**<prefix>.unitigs.gfa** Contigs split at bubble intersections.

**<prefix>.unitigs.bed** The position of each unitig in a contig.

## METADATA

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.

**<prefix>.contigs.layout, <prefix>.unitigs.layout** (undocumented)

**<prefix>.contigs.layout.readToTig, <prefix>.unitigs.layout.readToTig** The position of each read in a contig (unitig).

**<prefix>.contigs.layout.tigInfo, <prefix>.unitigs.layout.tigInfo** A list of the contigs (unitigs), lengths, coverage, number of reads and other metadata. Essentially the same information provided in the FASTA header line.



## CHAPTER 4

---

### Canu Pipeline

---

The pipeline is described in 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](#). bioRxiv. (2016). Figure 1 of the paper shows the primary pipeline (below, top) and the supplement contains the sub-pipeline for building read and overlap databases (below, bottom).



---

## Canu Parameter Reference

---

To get the most up-to-date options, run

```
canu -options
```

The default values below will vary based on the input data type and genome size.

Boolean options accept true/false or 1/0.

Memory sizes are assumed to be in gigabytes if no units are supplied. Values may be non-integer with or without a unit - 'k' for kilobytes, 'm' for megabytes, 'g' for gigabytes or 't' for terabytes. For example, "0.25t" is equivalent to "256g" (or simply "256").

### 5.1 Global Options

The catch all category.

**errorRate <float=unset> (OBSOLETE)** This parameter was removed on January 27th, 2016, and is valid only in Canu 1.4 or earlier. Canu currently still accepts the *errorRate* parameter, but its use is strongly discouraged.

The expected error in a single corrected read. The seven error rates were then set to three times this value (except for *corErrorRate*).

**rawErrorRate <float=unset>** The allowed difference in an overlap between two uncorrected reads, expressed as fraction error. Sets *corOvlErrorRate* and *corErrorRate*. The *rawErrorRate* typically does not need to be modified. It might need to be increased if very early reads are being assembled. The default is 0.300 For PacBio reads, and 0.500 for Nanopore reads.

**correctedErrorRate <float=unset>** The allowed difference in an overlap between two corrected reads, expressed as fraction error. Sets *obtOvlErrorRate*, *utgOvlErrorRate*, *obtErrorRate*, *utgErrorRate*, and *cnsErrorRate*. The *correctedErrorRate* can be adjusted to account for the quality of read correction, for the amount of divergence in the sample being assembled, and for the amount of sequence being assembled. The default is 0.045 for PacBio reads, and 0.144 for Nanopore reads.

For low coverage datasets (less than 30X), we recommend increasing *correctedErrorRate* slightly, by 1% or so.

For high-coverage datasets (more than 60X), we recommend decreasing *correctedErrorRate* slightly, by 1% or so.

Raising the *correctedErrorRate* will increase run time. Likewise, decreasing *correctedErrorRate* will decrease run time, at the risk of missing overlaps and fracturing the assembly.

**minReadLength <integer=1000>** Reads shorter than this are not loaded into the assembler. Reads output by correction and trimming that are shorter than this are discarded.

Must be no smaller than minOverlapLength.

If set high enough, the gatekeeper module will claim there are errors in the input reads, as too many of the input reads have been discarded. As long as there is sufficient coverage, this is not a problem. See *stopOnReadQuality* and *stopOnLowCoverage*

**minOverlapLength <integer=500>** Overlaps shorter than this will not be discovered. Smaller values can be used to overcome lack of read coverage, but will also lead to false overlaps and potential misassemblies. Larger values will result in more correct assemblies, but more fragmented, assemblies.

Must be no bigger than minReadLength.

**readSamplingCoverage <integer=unset>** After loading all reads into the sequence store, flag some reads as ‘not to be used’ until this amount of coverage remains. Reads are flagged according to the score described in *readSamplingBias*.

**readSamplingBias <float=0.0>** Adjust the sampling bias towards shorter (negative numbers) or longer (positive numbers) reads. Reads are assigned a score of  $random * length ^ bias$  and the lowest scoring reads are flagged as described in *readSamplingCoverage*.

**genomeSize <float=unset> *required*** An estimate of the size of the genome. Common suffices are allowed, for example, 3.7m or 2.8g.

The genome size estimate is used to decide how many reads to correct (via the *corOutCoverage* parameter) and how sensitive the mhap overlapper should be (via the *mhapSensitivity* parameter). It also impacts some logging, in particular, reports of NG50 sizes.

**canuIteration <internal parameter, do not use>** Which parallel iteration is being attempted.

**canuIterationMax <integer=2>** How many parallel iterations to try. Ideally, the parallel jobs, run under grid control, would all finish successfully on the first try. Sometimes, jobs fail due to other jobs exhausting resources (memory), or by the node itself failing. In this case, canu will launch the jobs again. This parameter controls how many times it tries.

**onSuccess <string=unset>** Execute the command supplied when Canu successfully completes an assembly. The command will execute in the <assembly-directory> (the -d option to canu) and will be supplied with the name of the assembly (the -p option to canu) as its first and only parameter.

**onFailure <string=unset>** Execute the command supplied when Canu terminates abnormally. The command will execute in the <assembly-directory> (the -d option to canu) and will be supplied with the name of the assembly (the -p option to canu) as its first and only parameter.

There are two exceptions when the command is not executed: if a ‘spec’ file cannot be read, or if canu tries to access an invalid parameter. The former will be reported as a command line error, and canu will never start. The latter should never occur except when developers are developing the software.

## 5.2 Process Control

**showNext <boolean=false>** Report the first major command that would be run, but don’t run it. Processing to get to that command, for example, checking the output of the previous command or preparing inputs for the next command, is still performed.



**stopOnReadQuality** <string=false> If set, Canu will stop with the following error if there are significantly fewer reads or bases loaded into the read store than what is in the input data.

```
Gatekeeper detected potential problems in your input reads.

Please review the logging in files:
  /assembly/godzilla/asm.gkpStore.BUILDING.err
  /assembly/godzilla/asm.gkpStore.BUILDING/errorLog

If you wish to proceed, rename the store with the following command and restart
↳ canu.

mv /assembly/godzilla/asm.gkpStore.BUILDING \
  /assembly/godzilla/asm.gkpStore.ACCEPTED

Option stopOnReadQuality=false skips these checks.
```

The missing reads could be too short (decrease *minReadLength* to include them), or have invalid bases or quality values. A summary of the files loaded and errors detected is in the `asm.gkpStore.BUILDING.err` file, with full gory details in the `asm.gkpStore.BUILDING/errorLog`.

To proceed, set `stopOnReadQuality=false` or rename the directory as shown.

Note that *U* bases are silently translated to *T* bases, to allow assembly of RNA sequences.

**stopOnLowCoverage** <integer=10> Stop the assembly if read coverage is too low to be useful. Coverage is checked when when input sequences are initially loaded into the sequence store, when corrected reads are generated, and when read ends are trimmed off.

**stopAfter** <string=undefined> If set, Canu will stop processing after a specific stage in the pipeline finishes.

Valid values for `stopAfter` are:

- `gatekeeper` - stops after the reads are loaded into the assembler read database.
- `meryl` - stops after frequent kmers are tabulated.
- `overlapConfigure` - stops after overlap jobs are configured.
- `overlap` - stops after overlaps are generated, before they are loaded into the overlap database.
- `overlapStoreConfigure` - stops after the jobs for creating the overlap store are configured.
- `overlapStore` - stops after overlaps are loaded into the overlap database.
- `readCorrection` - stops after corrected reads are generated.
- `readTrimming` - stops after trimmed reads are generated.
- `unitig` - stops after unitigs and contigs are created.
- `consensusConfigure` - stops after consensus jobs are configured.
- `consensus` - stops after consensus sequences are loaded into the databases.

## 5.3 General Options

**shell** <string="/bin/sh"> A path to a Bourne shell, to be used for executing scripts. By default, `/bin/sh`, which is typically the same as `bash`. C shells (`csh`, `tcsh`) are not supported.

**java** <string="java"> A path to a Java application launcher of at least version 1.8.

**gnuplot** <string="gnuplot"> A path to the gnuplot graphing utility. Plotting is disabled if this is unset (*gnuplot=* or *gnuplot=undef*), or if gnuplot fails to execute, or if gnuplot cannot generate plots. The latter two conditions generate warnings in the diagnostic output of Canu.

**gnuplotImageFormat** <string="png"> The type of image to generate in gnuplot. By default, canu will use png, svg or gif, in that order.

**preExec** <string=undef> A single command that will be run before Canu starts in a grid-enabled configuration. Can be used to set up the environment, e.g., with 'module'.

## 5.4 File Staging

The correction stage of Canu requires random access to all the reads. Performance is greatly improved if the gkpStore database of reads is copied locally to each node that computes corrected read consensus sequences. This 'staging' is enabled by supplying a path name to fast local storage with the stageDirectory option, and, optionally, requesting access to that resource from the grid with the gridEngineStageOption option.

**stageDirectory** <string=undefined> A path to a directory local to each compute node. The directory should use an environment variable specific to the grid engine to ensure that it is unique to each task.

For example, in Sun Grid Engine, */scratch/\$JOB\_ID-\$SGE\_TASK\_ID* will use both the numeric job ID and the numeric task ID. In SLURM, */scratch/\$SLURM\_JOBID* accomplishes the same.

If specified on the command line, be sure to escape the dollar sign, otherwise the shell will try to expand it before Canu sees the option: *stageDirectory=/scratch/\$JOB\_ID-\$SGE\_TASK\_ID*.

If specified in a specFile, do not escape the dollar signs.

**gridEngineStageOption** <string=undefined> This string is passed to the job submission command, and is expected to request local disk space on each node. It is highly grid specific. The string *DISK\_SPACE* will be replaced with the amount of disk space needed, in gigabytes.

On SLURM, an example is *-gres=lscratch:DISK\_SPACE*

## 5.5 Cleanup Options

**saveOverlaps** <boolean=false> If set to 'false', the raw overlapper outputs are removed as soon as they are loaded into an overlap store. Also, the correction and trimming overlap stores are removed when they are no longer needed.. This is recommended in nearly every case.

If set to 'stores', the raw overlapper outputs are removed, but all of the overlap stores are retained. The overlap stores capture all the critical information in the raw outputs and the raw outputs are redundant and unwieldy. Retaining the overlap stores can allow one to 'back up' and redo a step, but this is generally not useful unless one is familiar with the algorithms.

If set to 'true', all overlapper outputs and all stores are retained. This is useful for debugging potential problems with the overlap store.

**saveReadCorrections** <boolean=false>. If set, do not remove raw corrected read output from correction/2-correction. Normally, this output is removed once the corrected reads are generated.

**saveIntermediates** <boolean=false> If set, do not remove intermediate outputs. Normally, intermediate files are removed once they are no longer needed.

NOT IMPLEMENTED.

**saveMerCounts** <boolean=false> If set, do not remove meryl binary databases.

**saveReads** <boolean=false> If set, save the corrected reads (in `asm.correctedReads.fasta.gz`) and trimmed reads (in `asm.trimmedReads.fasta.gz`). Both read sets are saved in the `asm.gkpStore`, and can be retrieved later.

## 5.6 Executive Configuration

The Canu ‘executive’ is responsible for controlling what tasks run and when they run. It doesn’t directly do any significant computations, rather it just examines the files that exist and decides which component to run next. For example, if overlaps exist but contigs do not, it would create contigs next.

When under grid control, some tasks can be run in the same job as the executive, if there is enough memory and threads reserved for the executive. The benefit of this is slight; on a heavily loaded grid, it would reduce the number of job scheduling iterations Canu needs to run.

**executiveMemory** <integer=4>

The amount of memory, in gigabytes, to reserve when running the Canu executive (and any jobs it runs directly). Increasing this past 4 GB can allow some tasks (such as creating an overlap store or creating contigs) to run directly, without needing a separate grid job.

**executiveThreads** <integer=1>

The number of threads to reserve for the Canu executive.

## 5.7 Overlapper Configuration

Overlaps are generated for three purposes: read correction, read trimming and unitig construction. The algorithm and parameters used can be set independently for each set of overlaps.

Two overlap algorithms are in use. One, `mhap`, is typically applied to raw uncorrected reads and returns alignment-free overlaps with imprecise extents. The other, the original overlapper algorithm ‘`ovl`’, returns alignments but is much more expensive.

There are three sets of parameters, one for the ‘`mhap`’ algorithm, one for the ‘`ovl`’ algorithm, and one for the ‘`minimap`’ algorithm. Parameters used for a specific type of overlap are set by a prefix on the option: ‘`cor`’ for read correction, ‘`obt`’ for read trimming (‘overlap based trimming’) or ‘`utg`’ for unitig construction. For example, ‘`corOverlapper=ovl`’ would set the overlapper used for read correction to the ‘`ovl`’ algorithm.

**{prefix}Overlapper** <string=see-below> Specify which overlap algorithm, ‘`mhap`’ or ‘`ovl`’ or ‘`minimap`’. The default is to use ‘`mhap`’ for ‘`cor`’ and ‘`ovl`’ for both ‘`obt`’ and ‘`utg`’.

### 5.7.1 Overlapper Configuration, ovl Algorithm

**{prefix}OvlErrorRate** <float=unset> Overlaps above this error rate are not computed. \* *corOvlErrorRate* applies to overlaps generated for correcting reads; \* *obtOvlErrorRate* applied to overlaps generated for trimming reads; \* *utgOvlErrorRate* applies to overlaps generated for assembling reads. These limits apply to the ‘`ovl`’ overlap algorithm and when alignments are computed for `mhap` overlaps with *mhapReAlign*.

**{prefix}OvlFrequentMers** <string=undefined> Do not seed overlaps with these kmers (fasta format).

**{prefix}OvlHashBits** <integer=unset> Width of the kmer hash. Width 22=1gb, 23=2gb, 24=4gb, 25=8gb. Plus 10b per `ovlHashBlockLength`.

**{prefix}OvlHashBlockLength** <integer=unset> Amount of sequence (bp) to load into the overlap hash table.

**{prefix}OvlHashLoad** <integer=unset> Maximum hash table load. If set too high, table lookups are inefficient; if too low, search overhead dominates run time.

**{prefix}OvlMerDistinct <integer=unset>** K-mer frequency threshold; the least frequent fraction of distinct mers can seed overlaps.

**{prefix}OvlMerSize <integer=unset>** K-mer size for seeds in overlaps.

**{prefix}OvlMerThreshold <integer=unset>** K-mer frequency threshold; mers more frequent than this count are not used to seed overlaps.

**{prefix}OvlMerTotal <integer=unset>** K-mer frequency threshold; the least frequent fraction of all mers can seed overlaps.

**{prefix}OvlRefBlockLength <integer=unset>** Amount of sequence (bp to search against the hash table per batch.

**{prefix}OvlRefBlockSize <integer=unset>** Number of reads to search against the hash table per batch.

## 5.7.2 Overlapper Configuration, mhap Algorithm

**{prefix}MhapBlockSize <integer=unset>** For the MHAP overlapper, the number of reads to load per GB of memory (mhapMemory). When mhapSensitivity=high, this value is automatically divided by two.

**{prefix}MhapMerSize <integer=unset>** K-mer size for seeds in mhap.

**{prefix}ReAlign <boolean=false>** Compute actual alignments from mhap overlaps. uses either obtErrorRate or ovlErrorRate, depending on which overlaps are computed)

**{prefix}MhapSensitivity <string="normal">** Coarse sensitivity level: ‘low’, ‘normal’ or ‘high’. Based on read coverage (which is impacted by genomeSize), ‘low’ sensitivity is used if coverage is more than 60; ‘normal’ is used if coverage is between 60 and 30, and ‘high’ is used for coverages less than 30.

## 5.7.3 Overlapper Configuration, mmap Algorithm

**{prefix}MMapBlockSize <integer=unset>** Number of reads per 1GB block. Memory \* size is loaded into memory per job.

**{prefix}MMapMerSize <integer=unset>** K-mer size for seeds in minimap.

## 5.8 Overlap Store

The overlap algorithms return overlaps in an arbitrary order, however, all other algorithms (or nearly all) require all overlaps for a single read to be readily available. Thus, the overlap store collects and sorts the overlapper outputs into a store of overlaps, sorted by the first read in the overlap. Each overlap is listed twice in the store, once in an “A vs B” format, and once in a “B vs A” format (that is, swapping which read is ‘first’ in the overlap description).

Two construction algorithms are supported. A ‘sequential’ method uses a single data stream, and is faster for small and moderate size assemblies. A ‘parallel’ method uses multiple compute nodes and can be faster (depending on your network disk bandwidth) for moderate and large assemblies. Be advised that the parallel method is less efficient than the sequential method, and can easily thrash consumer-level NAS devices resulting in exceptionally poor performance.

The sequential method load all overlapper outputs (.ovb files in 1-overlapper) into memory, duplicating each overlap. It then sortes overlaps, and creates the final overlap store.

The parallel method uses two parallel tasks: bucketizing (‘ovb’ tasks) and sorting (‘ovs’ tasks). Bucketizing reads the outputs of the overlap tasks (ovb files in 1-overlapper), duplicates each overlap, and writes these to intermediate files. Sorting tasks load these intermediate file into memory, sorts the overlaps, then writes the sorted overlaps back to disk. There will be one ‘bucketizer’ (‘ovb’ tasks) task per overlap task, and tens to hundreds of ‘sorter’ (‘ovs’ tasks). A final ‘indexing’ step is done in the Canu executive, which ties all the various files together into the final overlap store.

Increasing `ovsMemory` will allow more overlaps to fit into memory at once. This will allow larger assemblies to use the sequential method, or reduce the number of ‘ovs’ tasks for the parallel method.

Increasing the allowed memory for the Canu executive can allow the overlap store to be constructed as part of the executive job – a separate grid job for constructing the store is not needed.

**ovsMemory** <float> How much memory, in gigabytes, to use for constructing overlap stores. Must be at least 256m or 0.25g.

## 5.9 Meryl

The ‘meryl’ algorithm counts the occurrences of kmers in the input reads. It outputs a FASTA format list of frequent kmers, and (optionally) a binary database of the counts for each kmer in the input.

Meryl can run in (almost) any memory size, by splitting the computation into smaller (or larger) chunks.

**merylMemory** <integer=unset> Amount of memory, in gigabytes, to use for counting kmers.

**merylThreads** <integer=unset> Number of compute threads to use for kmer counting.

## 5.10 Overlap Based Trimming

**obtErrorRate** <float=unset> Stringency of overlaps to use for trimming reads.

**trimReadsOverlap** <integer=1> Minimum overlap between evidence to make contiguous trim.

**trimReadsCoverage** <integer=1> Minimum depth of evidence to retain bases.

## 5.11 Grid Engine Support

Canu directly supports most common grid scheduling systems. Under normal use, Canu will query the system for grid support, configure itself for the machines available in the grid, then submit itself to the grid for execution. The Canu pipeline is a series of about a dozen steps that alternate between embarrassingly parallel computations (e.g., overlap computation) and sequential bookkeeping steps (e.g., checking if all overlap jobs finished). This is entirely managed by Canu.

Canu has first class support for the various schedulers derived from Sun Grid Engine (Univa, Son of Grid Engine) and the Simple Linux Utility for Resource Management (SLURM), meaning that the developers have direct access to these systems. Platform Computing’s Load Sharing Facility (LSF) and the various schedulers derived from the Portable Batch System (PBS, Torque and PBSPro) are supported as well, but without developer access bugs do creep in. As of Canu v1.5, support seems stable and working.

**useGrid** <boolean=true> Master control. If ‘false’, no algorithms will run under grid control. Does not change the value of the other useGrid options.

If ‘remote’, jobs are configured for grid execution, but not submitted. A message, with commands to launch the job, is reported and canu halts execution.

Note that the host used to run canu for ‘remote’ execution must know about the grid, that is, it must be able to submit jobs to the grid.

It is also possible to enable/disable grid support for individual algorithms with options such as *useGridBAT*, *useGridCNS*, et cetera. This has been useful in the (far) past to prevent certain algorithms, notably overlap error adjustment, from running too many jobs concurrently and thrashing disk. Recent storage systems seem to be able to handle the load better – computers have gotten faster quicker than genomes have gotten larger.

There are many options for configuring a new grid ('gridEngine\*') and for configuring how canu configures its computes to run under grid control ('gridOptions\*'). The grid engine to use is specified with the 'gridEngine' option.

**gridEngine** <string> Which grid engine to use. Auto-detected. Possible choices are 'sge', 'pbs', 'pbspro', 'lsf' or 'slurm'.

### 5.11.1 Grid Engine Configuration

There are many options to configure support for a new grid engine, and we don't describe them fully. If you feel the need to add support for a new engine, please contact us. That said, file `src/pipeline/canu/Defaults.pm` lists a whole slew of parameters that are used to build up grid commands, they all start with `gridEngine`. For each grid, these parameters are defined in the various `src/pipeline/Grid_*.pm` modules. The parameters are used in `src/pipeline/canu/Execution.pm`.

For SGE grids, two options are sometimes necessary to tell canu about peculiarities of your grid: `gridEngineThreadsOption` describes how to request multiple cores, and `gridEngineMemoryOption` describes how to request memory. Usually, canu can figure out how to do this, but sometimes it reports an error such as:

```
-- WARNING: Couldn't determine the SGE parallel environment to run multi-threaded_
↳codes.
Valid choices are (pick one and supply it to canu):
gridEngineThreadsOption="-pe make THREADS"
gridEngineThreadsOption="-pe make-dedicated THREADS"
gridEngineThreadsOption="-pe mpich-rr THREADS"
gridEngineThreadsOption="-pe openmpi-fill THREADS"
gridEngineThreadsOption="-pe smp THREADS"
gridEngineThreadsOption="-pe thread THREADS"
```

or:

```
-- WARNING: Couldn't determine the SGE resource to request memory.
Valid choices are (pick one and supply it to canu):
gridEngineMemoryOption="-l h_vmem=MEMORY"
gridEngineMemoryOption="-l mem_free=MEMORY"
```

If you get such a message, just add the appropriate line to your canu command line. Both options will replace the uppercase text (THREADS or MEMORY) with the value canu wants when the job is submitted. For `gridEngineMemoryOption`, any number of `-l` options can be supplied; we could use `gridEngineMemoryOption="-l h_vmem=MEMORY -l mem_free=MEMORY"` to request both `h_vmem` and `mem_free` memory.

### 5.11.2 Grid Options

To run on the grid, each stage needs to be configured - to tell the grid how many cores it will use and how much memory it needs. Some support for this is automatic (for example, `overlapInCore` and `mhap` know how to do this), others need to be manually configured. Yes, it's a problem, and yes, we want to fix it.

The `gridOptions*` parameters supply grid-specific options to the grid submission command.

**gridOptions** <string=unset> Grid submission command options applied to all grid jobs

**gridOptionsJobName** <string=unset> Grid submission command jobs name suffix

**gridOptionsBAT** <string=unset> Grid submission command options applied to unitig construction with the bogart algorithm

**gridOptionsGFA** <string=unset> Grid submission command options applied to gfa alignment and processing

**gridOptionsCNS** <string=unset> Grid submission command options applied to unitig consensus jobs

**gridOptionsCOR** <string=unset> Grid submission command options applied to read correction jobs

**gridOptionsExecutive** <string=unset> Grid submission command options applied to master script jobs

**gridOptionsOEA** <string=unset> Grid submission command options applied to overlap error adjustment jobs

**gridOptionsRED** <string=unset> Grid submission command options applied to read error detection jobs

**gridOptionsOVb** <string=unset> Grid submission command options applied to overlap store bucketizing jobs

**gridOptionsOVS** <string=unset> Grid submission command options applied to overlap store sorting jobs

**gridOptionsCORMHAP** <string=unset> Grid submission command options applied to mhap overlaps for correction jobs

**gridOptionsCOROVL** <string=unset> Grid submission command options applied to overlaps for correction jobs

**gridOptionsOBTMHAP** <string=unset> Grid submission command options applied to mhap overlaps for trimming jobs

**gridOptionsOBTOVL** <string=unset> Grid submission command options applied to overlaps for trimming jobs

**gridOptionsUTGMHAP** <string=unset> Grid submission command options applied to mhap overlaps for unitig construction jobs

**gridOptionsUTGOVL** <string=unset> Grid submission command options applied to overlaps for unitig construction jobs

## 5.12 Algorithm Selection

Several algorithmic components of canu can be disabled, based on the type of the reads being assembled, the type of processing desired, or the amount of compute resources available. Overlap

**enableOEA** <boolean=true> Do overlap error adjustment - comprises two steps: read error detection (RED and overlap error adjustment (OEA

### 5.12.1 Algorithm Execution Method

Canu has a fairly sophisticated (or complicated, depending on if it is working or not) method for dividing large computes, such as read overlapping and consensus, into many smaller pieces and then running those pieces on a grid or in parallel on the local machine. The size of each piece is generally determined by the amount of memory the task is allowed to use, and this memory size – actually a range of memory sizes – is set based on the genomeSize parameter, but can be set explicitly by the user. The same holds for the number of processors each task can use. For example, a genomeSize=5m would result in overlaps using between 4gb and 8gb of memory, and between 1 and 8 processors.

Given these requirements, Canu will pick a specific memory size and number of processors so that the maximum number of jobs will run at the same time. In the overlapper example, if we are running on a machine with 32gb memory and 8 processors, it is not possible to run 8 concurrent jobs that each require 8gb memory, but it is possible to run 4 concurrent jobs each using 6gb memory and 2 processors.

To completely specify how Canu runs algorithms, one needs to specify a maximum memory size, a maximum number of processors, and how many pieces to run at one time. Users can set these manually through the {prefix}Memory, {prefix}Threads and {prefix}Concurrency options. If they are not set, defaults are chosen based on genomeSize.

**{prefix}Concurrency** <integer=unset> Set the number of tasks that can run at the same time, when running without grid support.

**{prefix}Threads <integer=unset>** Set the number of compute threads used per task.

**{prefix}Memory <integer=unset>** Set the amount of memory, in gigabytes, to use for each job in a task.

Available prefixes are:

Prefix		Algorithm
cor obt utg	mhap	Overlap generation using the 'mhap' algorithm for 'cor'=correction,, 'obt'=trimming or 'utg'=assembly.
cor obt utg	mmap	Overlap generation using the 'minimap' algorithm for 'cor'=correction,, 'obt'=trimming or 'utg'=assembly.
cor obt utg	ovl	Overlap generation using the 'overlapInCore' algorithm for 'cor'=correction,, 'obt'=trimming or 'utg'=assembly.
	ovb	Parallel overlap store bucketizing
	ovs	Parallel overlap store bucket sorting
	cor	Read correction
	red	Error detection in reads
	oea	Error adjustment in overlaps
	bat	Unitig/contig construction
	cns	Unitig/contig consensus

For example, 'mhapMemory' would set the memory limit for computing overlaps with the mhap algorithm; 'cormhap-Memory' would set the memory limit only when mhap is used for generating overlaps used for correction.

The 'minMemory', 'maxMemory', 'minThreads' and 'maxThreads' options will apply to all jobs, and can be used to artificially limit canu to a portion of the current machine. In the overlapper example above, setting maxThreads=4 would result in two concurrent jobs instead of four.

## 5.13 Overlap Error Adjustment

red = Read Error Detection oea = Overlap Error Adjustment

**oeaBatchLength <unset>** Number of bases per overlap error correction batch

**oeaBatchSize <unset>** Number of reads per overlap error correction batch

**redBatchLength <unset>** Number of bases per fragment error detection batch

**redBatchSize <unset>** Number of reads per fragment error detection batch



## 5.14 Unitigger

**unitigger** <string="bogat"> Which unitig construction algorithm to use. Only "bogat" is supported.

**utgErrorRate** <float=unset> Stringency of overlaps used for constructing contigs. The *bogat* algorithm uses the distribution of overlap error rates to filter high error overlaps; *bogat* will never see overlaps with error higher than this parameter.

**batOptions** <unset> Advanced options to bogart

## 5.15 Consensus Partitioning

STILL DONE BY UNITIGGER, NEED TO MOVE OUTSIDE

**cnsConsensus** Which algorithm to use for computing consensus sequences. Only 'utgcns' is supported.

**cnsPartitions** Compute consensus by splitting the tigs into N partitions.

**cnsPartitionMin** Don't make a partition with fewer than N reads

**cnsMaxCoverage** Limit unitig consensus to at most this coverage.

**cnsErrorRate** Inform the consensus generation algorithm of the amount of difference it should expect in a read-to-read alignment. Typically set to *utgOvlErrorRate*. If set too high, reads could be placed in an incorrect location, leading to errors in the consensus sequence. If set too low, reads could be omitted from the consensus graph (or multialignment, depending on algorithm), resulting in truncated consensus sequences.

## 5.16 Read Correction

The first step in Canu is to find high-error overlaps and generate corrected sequences for subsequent assembly. This is currently the fastest step in Canu. By default, only the longest 40X of data (based on the specified genome size) is used for correction. Typically, some reads are trimmed during correction due to being chimeric or having erroneous sequence, resulting in a loss of 20-25% (30X output). You can force correction to be non-lossy by setting *corMinCoverage=0*, in which case the corrected reads output will be the same length as the input data, keeping any high-error unsupported bases. Canu will trim these in downstream steps before assembly.

If you have a dataset with uneven coverage or small plasmids, correcting the longest 40X may not give you sufficient coverage of your genome/plasmid. In these cases, you can set *corOutCoverage=999*, or any value greater than your total input coverage which will correct and assemble all input data, at the expense of runtime.

**corErrorRate** <integer=unset> Do not use overlaps with error rate higher than this (estimated error rate for *mhap* and *minimap* overlaps).

**corConsensus** <string="falconpipe"> Which algorithm to use for computing read consensus sequences. Only 'falcon' and 'falconpipe' are supported.

**corPartitions** <integer=128> Partition read correction into N jobs

**corPartitionMin** <integer=25000> Don't make a read correction partition with fewer than N reads

**corMinEvidenceLength** <integer=unset> Limit read correction to only overlaps longer than this; default: unlimited

**corMinCoverage** <integer=4> Limit read correction to regions with at least this minimum coverage. Split reads when coverage drops below threshold.

**corMaxEvidenceError** <integer=unset> Limit read correction to only overlaps at or below this fraction error; default: unlimited

**corMaxEvidenceCoverageGlobal** <string="1.0x"> Limit reads used for correction to supporting at most this coverage; default: 1.0 \* estimated coverage

**corMaxEvidenceCoverageLocal** <string="2.0x"> Limit reads being corrected to at most this much evidence coverage; default: 10 \* estimated coverage

**corOutCoverage** <integer=40> Only correct the longest reads up to this coverage; default 40

**corFilter** <string="expensive"> Method to filter short reads from correction; 'quick' or 'expensive' or 'none'

## 5.17 Output Filtering

**contigFilter** <minReads, integer=2> <minLength, integer=0> <singleReadSpan, float=1.0> <lowCovSpan, float=0.5> <lowCovDepth, integer=5>

A contig that needs any of the following conditions is flagged as 'unassembled' and removed from further consideration:

- fewer than minReads reads (default 2)
- shorter than minLength bases (default 0)
- a single read covers more than singleReadSpan fraction of the contig (default 1.0)
- more than lowCovSpan fraction of the contig is at coverage below lowCovDepth (defaults 0.5, 5)

This filtering is done immediately after initial contigs are formed, before potentially incorrectly spanned repeats are detected. Initial contigs that incorrectly span a repeat can be split into multiple contigs; none of these new contigs will be flagged as 'unassembled', even if they are a single read.

---

## Canu Command Reference

---

Every command, even the useless ones.

Commands marked as ‘just usage’ were automatically generated from the command line usage summary. Yes, some of them even crashed.

**commands/bogart (just usage)** The unitig construction algorithm. BOG stands for Best Overlap Graph; we haven’t figured out what ART stands for.

**commands/bogus (just usage)** A unitig construction algorithm simulator. Given reads mapped to a reference, returns the largest unitigs possible.

**commands/canu (just usage)** The executive in charge! Coordinates all these commands to make an assembler.

**commands/correctOverlaps (just usage)** Part of Overlap Error Adjustment, recomputes overlaps given a set of read corrections.

**commands/estimate-mer-threshold (just usage)** Decides on a k-mer threshold for overlapInCore seeds.

**commands/fastqAnalyze (just usage)** Analyzes a FASTQ file and reports the best guess of the QV encoding. Can also rewrite the FASTQ to be in Sanger QV format.

**commands/fastqSample (just usage)** Extracts random reads from a single or mated FASTQ file. Extracts based on desired coverage, desired number of reads/pairs, desired fraction of the total, or desired total length.

**commands/fastqSimulate (just usage)** Creates reads with unbiased errors from a FASTA sequence.

**commands/fastqSimulate-sort (just usage)** Given input from fastqSimulate, sorts the reads by position in the reference.

**commands/filterCorrectionOverlaps (just usage)** Part of Read Correction, filters overlaps that shouldn’t be used for correcting reads.

**commands/findErrors (just usage)** Part of Overlap Error Adjustment, generates a multialignment for each read, outputs a list of suspected errors in the read.

**commands/gatekeeperCreate (just usage)** Loads FASTA or FASTQ reads into the canu read database, gkpStore.

**commands/gatekeeperDumpFASTQ (just usage)** Outputs FASTQ reads from the canu read database, gkpStore.

- commands/gatekeeperDumpMetaData (just usage)** Outputs read and library metadata from the canu read database, gkpStore.
- commands/gatekeeperPartition (just usage)** Part of Consensus, rearranges the canu read database, gkpStore, to localize read to unitigs.
- commands/generateCorrectionLayouts (just usage)** Part of Read Correction, generates the multialignment layout used to correct reads.
- commands/leaff (just usage)** Not actually part of canu, but it came along with meryl. Provides random access to FASTA, FASTQ and gkpStore. Also does some analysis tasks. Handy Swiss Army knife type of tool.
- commands/meryl (just usage)** Counts k-mer occurrences in FASTA, FASTQ and gkpStore. Performs mathematical and logical operations on the resulting k-mer databases.
- commands/mhapConvert (just usage)** Convert mhap output to overlap output.
- commands/ovStoreBucketizer (just usage)** Part of the parallel overlap store building pipeline, loads raw overlaps from overlapper into the store.
- commands/ovStoreBuild (just usage)** Sequentially builds an overlap store from raw overlaps. Simplest to run, but slow on large datasets.
- commands/ovStoreDump (just usage)** Dumps overlaps from the overlap store, ovlStore.
- commands/ovStoreIndexer (just usage)** Part of the parallel overlap store building pipeline, finalizes the store, after sorting with ovStoreSorter.
- commands/ovStoreSorter (just usage)** Part of the parallel overlap store building pipeline, sorts overlaps loaded into the store by ovStoreBucketizer.
- commands/overlapConvert (just usage)** Reads raw overlapper output, writes overlaps as ASCII. The reverse of overlapImport.
- commands/overlapImport (just usage)** Reads ASCII overlaps in a few different formats, writes either ‘raw overlapper output’ or creates an ovlStore.
- commands/overlapInCore (just usage)** The classic overlapper algorithm.
- commands/overlapInCorePartition (just usage)** Generate partitioning to run overlapInCore jobs in parallel.
- commands/overlapPair (just usage)** An *experimental* algorithm to recompute overlaps and output the alignments.
- commands/prefixEditDistance-matchLimitGenerate (just usage)** Generate source code files with data representing the minimum length of a good overlap given some number of errors.
- commands/splitReads (just usage)** Part of Overlap Based Trimming, splits reads based on overlaps, specifically, looking for PacBio hairpin adapter signatures.
- commands/tgStoreCoverageStat (just usage)** Analyzes tigs in the tigStore, computes the classic [arrival rate statistic](#).
- commands/tgStoreDump (just usage)** Analyzes and outputs tigs from the tigStore, in various formats (FASTQ, layouts, multialignments, etc).
- commands/tgStoreFilter (just usage)** Analyzes tigs in the tigStore, marks those that appear to be spurious ‘degenerate’ tigs.
- commands/tgStoreLoad (just usage)** Loads tigs into a tigStore.
- commands/tgTigDisplay (just usage)** Displays the tig contained in a binary multialignment file, as output by utgcn.
- commands/trimReads (just usage)** Part of Overlap Based Trimming, trims reads based on overlaps.
- commands/utgcn (just usage)** Generates a multialignment for a tig, based on the layout stored in tigStore. Outputs FASTQ, layouts and binary multialignment files.

---

## Software Background

---

Canu is derived from [Celera Assembler](#), which is no longer maintained.

Celera Assembler [Myers 2000] was designed to reconstruct mammalian chromosomal DNA sequences from the short fragments of a whole genome shotgun sequencing project. Celera Assembler was used to produce reconstructions of several large genomes, namely those of *Homo sapiens* [Venter 2001], *Mus musculus* [Mural 2002], *Rattus norvegicus* [unpublished data], *Canis familiaris* [Kirkness 2003], *Drosophila melanogaster* [Adams 2000], and *Anopheles gambiae* [Holt 2001]. Celera Assembler was shown to be very accurate when its reconstruction of the human genome was compared to independent reconstructions completed later [Istrail 2004]. It was used to reconstructing one of the first large-scale metagenomic projects [Venter 2004, Rusch 2007] and a diploid human reference [Levy 2007, Denisov 2008]. It was adapted to 454 Pyrosequencing [Miller 2008] and PacBio sequencing [Koren 2012], demonstrating finished bacterial genomes [Koren 2013] and efficient algorithms for eukaryotic assembly [Berlin 2015].

Celera Assembler was released under the GNU General Public License, version 2 as a supplement to [Istrail 2004].

Canu [Koren and Walenz 2017] was branched from Celera Assembler in 2015, and specialized for single-molecule high-noise sequences. For the most recent license information on Canu, see [README.licences](#).

## 7.1 References

- Adams et al. (2000) *The Genome Sequence of Drosophila melanogaster*. *Science* 287 2185-2195.
- Myers et al. (2000) *A Whole-Genome Assembly of Drosophila*. *Science* 287 2196-2204.
- Venter et al. (2001) *The Sequence of the Human Genome*. *Science* 291 1304-1351.
- Mural et al. (2002) *A Comparison of Whole-Genome Shotgun-Derived Mouse Chromosome 16 and the Human Genome*. *Science* 296 1661-1671.
- Holt et al. (2002) *The Genome Sequence of the Malaria Mosquito Anopheles gambiae*. *Science* 298 129-149.
- Istrail et al. (2004) *Whole Genome Shotgun Assembly and Comparison of Human Genome Assemblies*. *PNAS* 101 1916-1921.
- Kirkness et al. (2003) *The Dog Genome: Survey Sequencing and Comparative Analysis*. *Science* 301 1898-1903.

- Venter et al. (2004) Environmental genome shotgun sequencing of the Sargasso Sea. *Science* 304 66-74.
- Levy et al. (2007) The Diploid Genome Sequence of an Individual Human. *PLoS Biology* 0050254
- Rusch et al. (2007) The Sorcerer II Global Ocean Sampling Expedition: Northwest Atlantic through Eastern Tropical Pacific. *PLoS Biology* 1821060.
- Denisov et al. (2008) Consensus Generation and Variant Detection by Celera Assembler. *Bioinformatics* 24(8):1035-40
- Miller et al. (2008) Aggressive Assembly of Pyrosequencing Reads with Mates. *Bioinformatics* 24(24):2818-2824
- Koren et al. (2012) Hybrid error correction and de novo assembly of single-molecule sequencing reads. *Nature Biotechnology*, July 2012.
- Koren et al. (2013) Reducing assembly complexity of microbial genomes with single-molecule sequencing. *Genome Biology* 14:R101.
- Berlin et. al. (2015) Assembling Large Genomes with Single-Molecule Sequencing and Locality Sensitive Hashing. *Nature Biotechnology*. (2015).
- Koren and Walenz et al. (2017) Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Research*. (2017).

Canu is a fork of the Celera Assembler designed for high-noise single-molecule sequencing (such as the PacBio RSII or Oxford Nanopore MinION).

## CHAPTER 8

---

### Publication

---

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*. (2017).





## CHAPTER 9

---

### Install

---

The easiest way to get started is to download a [release](#). If you encounter any issues, please report them using the [github issues page](#).

Alternatively, you can also build the latest unreleased from github:

```
git clone https://github.com/marbl/canu.git
cd canu/src
make -j <number of threads>
```



## CHAPTER 10

---

### Learn

---

- *Quick Start* - no experience or data required, download and assemble *Escherichia coli* today!
- *FAQ* - Frequently asked questions
- *Canu tutorial* - a gentle introduction to the complexities of canu.
- *Canu pipeline* - what, exactly, is canu doing, anyway?
- *Canu Parameter Reference* - all the parameters, grouped by function.
- *Canu Command Reference* - all the commands that canu runs for you.
- *Canu History* - the history of the Canu project.