# Genome assembly

## Genome assembly

Genome assembly is the process of converting short reads into a detailed set of sequences corresponding to the chromosome(s) of an organism.

To learn more about assembly visit
**http://www.ncbi.nlm.nih.gov/assembly/**
**http://www.ncbi.nlm.nih.gov/assembly/basics/**

**Assembly**

Genome assembly organization and additional information.

| Using Assembly | Submitting an Assembly | Related Resources |
|---|---|---|
| Assembly Help | Submission Information | Genome |
| Browse by Organism | Submission FAQ | Genome Reference Consortium |
| NCBI Assembly Data Model | AGP Specifications | Genome Remapping Service (Remap) |
| Assembly Basics | AGP Validation | |
| Genomes Download FAQ | | |
| Genomes FTP Site | | |

## Genome assembly: relevance

- Genome assembly is needed when a genome is first sequenced. We can relate reads to chromosomes.

- For the human genome, the assembly is "frozen" as a snapshot every few years. The current assembly is GRCh38. (GRC refers to Genome Reference Consortium at http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/)

- For most human genome work we do not need to do "de novo" (from anew) assembly. Instead we map reads to a reference genome—one that is already assembled.

- Genome assembly is a crucial behind-the-scenes part of calling human genome (or other) variants.

---

Whereas early genome assembly projects were often aided by clone maps or other mapping data, many current assembly projects forego these scaffolding data and only assemble genomes into smaller segments. Recently, new technologies have been invented that allow chromosome-scale assembly at a lower cost and faster speed than traditional methods.

Many new technologies can now be used to create chromosome-scale assemblies without costly and time-consuming methods such as BAC-end sequencing and physical mapping.

Rice and Green, 2019. New Approaches for Genome Assembly and Scaffolding

---

Consequently, the contiguity of new genome assemblies decreased as high-throughput sequencing was widely adopted (Figure 1b,c)
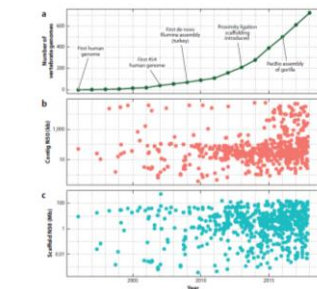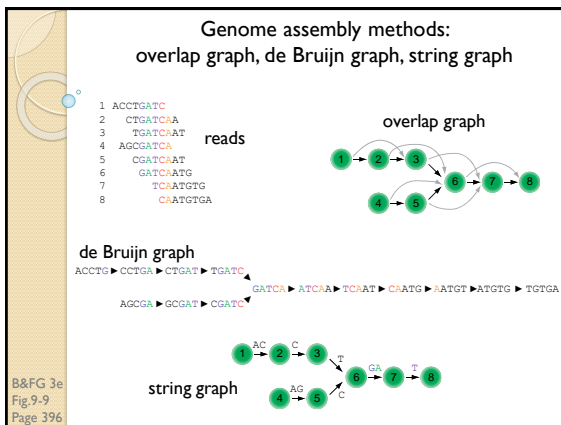
**Genome Contig Assembly**

No technology currently exists that can read DNA from one end to the other of even moderately sized chromosomes,which are typically tens or hundreds of millions of base pairs long. All current approaches for genome assembly read many segments that are considerably shorter than chromosomes.

Both long-read sequencing technologies implement single-molecule sequencing methods and generate reads with a distribution of lengths that, for assembly purposes, target a range of tens to hundreds of kilobases (kb)- typically 10–25 kb for PacBio HiFi reads (also circular consensus sequencing, CCS), 10–40 kb for PacBio continuous long reads (CLR) and 10 kb-2 Megabases (Mb) for ONT, where the upper limit is constrained principally by properties of the input material (Payne et al., 2018).

---

## Assembly algorithms

- overlap-layout-consensus: input DNA sequence reads are compared, all versus all, in the overlap step. The overlap-layout-consensus algorithm is based on identifying overlapping regions between reads and using these overlaps to construct longer contiguous sequences (contigs).
- de Bruijn graph: short words (k-mers) that are observed in the reads are the nodes of the graph, and edges are added when these k-mers are adjacent in sequence reads. In this process, each read is used to populate the graph but not compared directly to all the other reads.
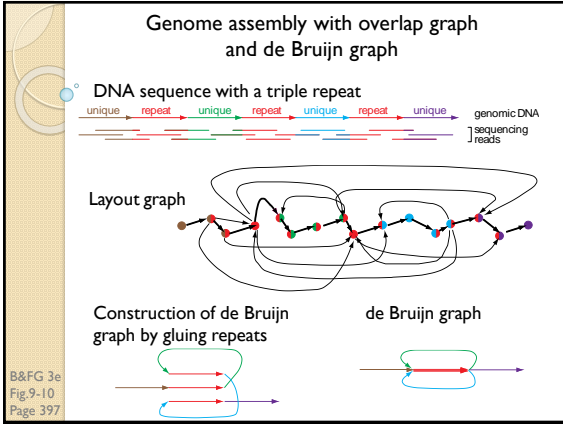- hybrid assembly

---

### Genome assembly methods:
### overlap graph, de Bruijn graph, string graph



B&FG 3e
Fig.9-9
Page 396

## Genome assembly with overlap graph and de Bruijn graph

**DNA sequence with a triple repeat**



**Layout graph**

**Construction of de Bruijn graph by gluing repeats**

**de Bruijn graph**

B&FG 3e
Fig.9-10
Page 397

---

Table 1 Commonly used assembly software

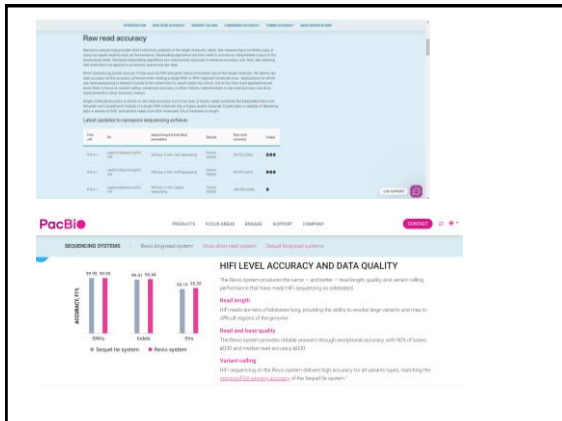| Software | URL and reference | Description |
|---|---|---|
| **Short-read assembly software** | | |
| Velvet | http://github.com/dzerbino/velvet (168) | Original de Bruijn graph assembler |
| SOAPdenovo | http://soap.genomics.org.cn/ (169) | De Bruijn graph assembler with error-correction step |
| Meraculous | https://jgi.doe.gov/data-and-tools/meraculous/ (170) | Hybrid k-mer/read-based |
| ALLPATHS-LG | http://software.broadinstitute.org/allpaths-lg/blog/ (171) | Uses unipath graph to collapse repeats |
| SGA | https://github.com/jts/sga (172) | Uses string graphs |
| AbySS | https://github.com/bcgsc/abyss (173) | Represents de Bruijn graph with a Bloom filter |
| DISCOVAR de novo | https://software.broadinstitute.org/software/discovar/blog/ (174) | Requires 250-bp PCR-free reads |
| Supernova | https://github.com/10XGenomics/supernova (149) | Assembles 10× linked reads |
| **Long-read assembly software** | | |
| HGAP | https://github.com/PacificBiosciences/Bioinformatics-Training/wiki/HGAP (124) | Error correction, overlap-layout-consensus assembly, and polishing workflow |
| Canu | https://github.com/marbl/canu (125) | K-mer-based overlap computation |
| FALCON | https://github.com/PacificBiosciences/FALCON (103) | Assembles phased diploid genomes |
| Flye | https://github.com/fenderglass/Flye (129) | Uses A-Bruijn graph |
| Minimm | https://github.com/lh3/minimm (126) | Fast, but no error correction |
| **Polishing software** | | |
| Pilon | https://github.com/broadinstitute/pilon (133) | Uses short-read alignments to correct errors |
| Arrow | https://github.com/PacificBiosciences/GenomicConsensus | Hidden Markov model and long-read alignments |
| Nanopolish | https://github.com/jts/nanopolish (115) | Nanopore only; uses original voltage data to correct errors |

Spades??

---

## CREATING MORE CONTIGUOUS ASSEMBLIES WITH LONG READS

- Pacific Biosciences (SMRT, 2009)
  - The incorporation of fluorescently labeled nucleotides is detected and reveals the sequence of the analysed DNA strand.
  - PacBio offers Continuous Long Reads (CLR) and Circular Consensus Sequencing (CCS) reads also called High-Fidelity (HiFi).
- Oxford Nanopore Sequencing (2005, 1 channel flow cell, etc.)
  - It works by monitoring changes to an electrical current as nucleic acids are passed through a protein nanopore. The resulting signal is decoded to provide the specific DNA or RNA sequence.

## Hybrid assembly

- The accuracy of the short reads is used to decrease the error rate of the long reads from up to 20% to as low as 0.1%. Then, the corrected long reads are assembled using an algorithm such as overlap-layout-consensus.

Is it still necessary with new chemistry used by ONT and PacBio?





https://gigabytejournal.com/articles/122

## NEW APPROACHES FOR LONG-RANGE GENOME SCAFFOLDING

- method called Hi-C, Omni-C (Hi-C is a chromosome conformation capture (3C)-based technology to detect pair-wise chromatin interactions genome-wide)
- Linked-Read Sequencing (single-tube long fragment reads (stLFR) and haplo-tagging (Meier et al., 2020; Wang et al., 2019)
- Optical maps
- Synteny-Based Methods



**Figure 4**

Overview of methods for long-range scaffolding. (a) In proximity ligation, chromatin is crosslinked and then restriction digested, ligated, and fragmented to create reads containing sequence from two different parts of the same chromosome. (b) In 10× linked-read sequencing, high-molecular weight DNA is combined with barcoded beads in oil droplets and then undergoes barcoding and amplification inside the droplets, resulting in reads with the same barcode that came from the same initial fragment of DNA. (c) BioNano optical maps are created by nicking high-molecular weight DNA with multiple nicking enzymes and attaching fluorescent markers at the nick sites. Contigs can then be aligned to the optical map by lining up nicking sequences in the contig with the locations of fluorescent markers in the map. (d) In synteny-based approaches, contigs are mapped to the assembled genomes of one or more related species. These alignments imply the order and orientation of the aligned contigs.

As of April 2021, four biochemical companies (Arima Genomics, Dovetail Genomics, Phase Genomics, and Qiagen) manufacture Hi-C kits, which are formulated with different components and protocols. In general, conventional Hi-C kits employ a restriction enzyme or a cocktail of multiple restriction enzymes, whereas Omni-C employs a sequence-independent endonuclease (Table 1). In Omni-C, to capture more proximal contacts, disuccinimidyl glutarate (DSG) and formaldehyde are used for sample fixation (Nowak et al., 2005), which is now provided as a kit by Dovetail Genomics.

https://onlinelibrary.wiley.com/doi/full/10.1111/mec.16146

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2858594/



https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6926122/



https://onlinelibrary.wiley.com/doi/full/10.1111/1755-0998.13312

Key approaches for genome assembly that are generally recommended in all species include the following:

(a) Genome assemblies should include long-read sequencing except in rare cases where it is effectively impossible to acquire adequately preserved samples needed for HMW DNA standards.

(b) At least one scaffolding approach should be included with genome assembly such as Hi-C mapping or optical mapping (linked-read data is also appropriate but may not be available for most future projects).

(c) Short-read data should be included for genome polishing, error correction, k-mer analyses, and estimating the percent of reads that map back to assembly.

https://onlinelibrary.wiley.com/doi/full/10.1111/1755-0998.13312

## Validation of genome assembly

- BUSCO (Benchmarking Universal SingleCopy Orthologs ) OrthoDB (BUSCO uses set of genes which are present in 90 % of species in one copy only)
- QUAST

## Genome annotation

- Genome annotation is the process of identifying and labeling functional elements within the genome, such as genes, regulatory regions, and repetitive elements.

**Table 3. Commonly used genome annotation tools and programs.**



Jung H, Ventura T, Chung JS, Kim WJ, Nam BH, et al. (2020) Twelve quick steps for genome assembly and annotation in the classroom. PLOS Computational Biology 16(11): e1008325. https://doi.org/10.1371/journal.pcbi.1008325
https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1008325

**Table 3. Commonly used genome annotation tools and programs.**



Jung H, Ventura T, Chung JS, Kim WJ, Nam BH, et al. (2020) Twelve quick steps for genome assembly and annotation in the classroom. PLOS Computational Biology 16(11): e1008325. https://doi.org/10.1371/journal.pcbi.1008325
https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1008325

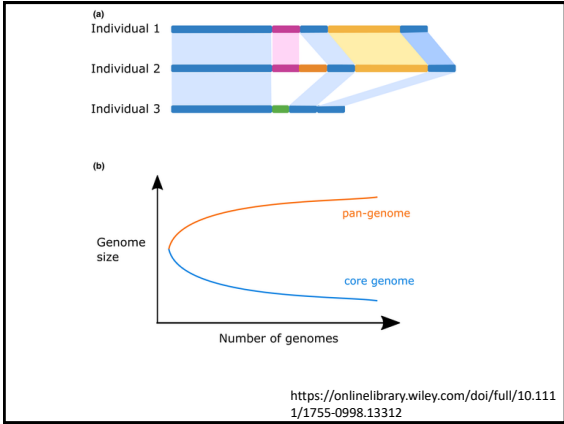**Table 3. Commonly used genome annotation tools and programs.**



Jung H, Ventura T, Chung JS, Kim WJ, Nam BH, et al. (2020) Twelve quick steps for genome assembly and annotation in the classroom. PLOS Computational Biology 16(11): e1008325. https://doi.org/10.1371/journal.pcbi.1008325
https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1008325

NCBI eukaryotic annotation pipeline

2015

2024

Conceptually similar to Ensembl pipeline

B&FG 3e
Fig. 15.14
Page 740

Biological challenges

- Repetitive regions (expanded gene families, complex repeats, highly repetitive regions such as centromeres and telomeres, and sex chromosomes, or at least portions of them.)
- Ploidy
- Pan and core genomes (The pan-genome represents all sequences among all of the DNA sequences that occur in a species whereas the core-genome is the DNA that is shared among all sequenced individuals.)
- *For example, the comparison of eight chromosome-level assemblies of Arabidopsis thaliana accessions revealed a core-genome, shared by all accessions, of ~105 Mb and ~24,000 genes, whereas the pan-genome was ~135 Mb in length and included ~30,000 genes (Jiao & Schneeberger, 2020), highlighting the vast amount of sequence data, including genes, that are missed by a single reference genome assembly.*

https://onlinelibrary.wiley.com/doi/full/10.1111/1755-0998.13312



(a)

Individual 1
Individual 2
Individual 3

(b)

Genome size

pan-genome

core genome

Number of genomes

https://onlinelibrary.wiley.com/doi/full/10.1111/1755-0998.13312

A    NGS variant calling

Single nucleotide variants (SNVs)

Small insertion/deletions (InDels)

Collapsed repeats

Inaccessible regions

B    long read variant calling

Deletion

Insertion

Tandem duplication

Inversion

C    *de novo* assembly

Chromosomal rearrangements
Accession1

Accession2

Separated haplophases
Haplo1

Haplo2

https://www.cambridge.org/core/journals/quantitative-plant-biology/article/plant-genome-sequence-assembly-in-the-era-of-long-reads-progress-challenges-and-future-directions/3932B956CC468D07D3E14525FEAFEC3E

## Why Is Chromosome-Scale Assembly Important?

- Cis-regulatory elements and the complexity of regulatory architecture
- Recombination
- Genetic association studies
- Chromosome evolution