## Austra & Lian Journal of Basic Sciences



australiansciencejournals.com/aljbs

E-ISSN: 2643-251X

**VOL 05 ISSUE 06 2024** 

# **Advances in Bioinformatics Tools for Genome Assembly**

#### Dr. Emily Zhang

Department of Computational Biology, University of Toronto, Canada

Email: emily.zhang@utoronto.ca

Abstract: The rapid advancement of next-generation sequencing (NGS) technologies has necessitated the parallel evolution of bioinformatics tools for efficient and accurate genome assembly. These tools are pivotal for reconstructing genomic sequences from short reads, enabling significant progress in genomics, personalized medicine, and evolutionary studies. This article reviews the recent advances in genome assembly tools, focusing on algorithmic innovations, hybrid assembly approaches, error correction techniques, and scalability for large genomes. It also highlights the integration of machine learning to enhance accuracy and discusses the future challenges in the field. The comprehensive overview aims to support researchers in selecting appropriate tools and methodologies for specific genome assembly projects.

**Keywords:** genome assembly, bioinformatics tools, next-generation sequencing, hybrid assembly, computational genomics

#### **INTRODUCTION:**

Genome assembly is a cornerstone of modern genomics, transforming raw sequence data into meaningful biological information. As sequencing platforms generate vast amounts of data with increasing speed and reduced cost, the need for efficient and accurate assembly tools has become more critical. Initially relying on Sanger sequencing, genome assembly methods have evolved to accommodate high-throughput platforms such as Illumina, PacBio, and Oxford Nanopore. However, the complexity of genomes, including repetitive elements and structural variations, presents considerable challenges. Recent advances in bioinformatics have led to the development of powerful assembly tools and algorithms capable of resolving such complexities. This paper explores key developments in genome assembly tools, examining their methodologies, applications, and future trajectories.

#### 1. Historical Evolution of Genome Assembly Tools:

Genome assembly tools have undergone several generations of development, each shaped by the prevailing sequencing technologies and computational limitations of the time. The objective of genome assembly is to reconstruct the original DNA sequence from fragments generated by sequencing platforms. However,

the path to this goal has not been linear — it has involved major paradigm shifts in both algorithmic frameworks and biological data characteristics.

## Stage 1: Overlap-Layout-Consensus (OLC) Models - The Sanger Era

The OLC model originated during the dominance of **Sanger sequencing**, which provided relatively **long** (~800–1000 bp), high-fidelity reads, but in **low throughput**. Genome assembly at this stage was akin to solving a jigsaw puzzle with fewer but larger and clearer pieces. OLC assemblers like **Celera Assembler**, which was famously used to assemble the human genome in 2001, worked by:

Identifying overlaps between all pairs of reads (using dynamic programming or hash-based methods),

Constructing a graph layout, where nodes represent reads and edges represent significant overlaps,

Generating a consensus sequence by aligning overlapping reads and resolving base discrepancies.

While effective for smaller genomes, OLC algorithms are **quadratic in time complexity**, making them computationally expensive for large datasets, especially when the number of reads exceeds millions. Additionally, **repeat regions**, which are prevalent in eukaryotic genomes, often confound layout steps by creating ambiguous paths.

## Stage 2: De Bruijn Graph Assemblers - The NGS Revolution

With the rise of **Next-Generation Sequencing (NGS)** technologies, especially **Illumina**, the sequencing landscape shifted toward generating **hundreds of millions to billions of short reads** (50–300 bp). These short reads are economical and have low error rates, but the sheer volume rendered OLC impractical due to scalability bottlenecks. This led to the adoption of **De Bruijn graph (DBG) methods**, which break sequences into overlapping k-mers (substrings of length k). In DBG:

Each k-mer is a node, and an edge is drawn from one k-mer to the next if they overlap by (k-1) bases.

This graph represents all possible paths through the data, and **Eulerian path algorithms** are used to reconstruct the genome.

Tools like **Velvet**, **SOAPdenovo**, **ABySS**, and later **SPAdes** utilized DBG efficiently. These methods are **linear in the number of k-mers**, allowing assemblies of large genomes on modest computational infrastructure. However, drawbacks include:

Choice of k being critical (too small: high connectivity; too large: low coverage),

Difficulty resolving complex repeats due to read length limitations,

**Sequencing errors** producing spurious nodes and branches.

To mitigate errors, these assemblers include **error correction algorithms** and **graph simplification techniques**, such as tip removal, bubble popping, and repeat resolution heuristics.

#### Stage 3: Long-Read Assemblers – Toward Greater Continuity

The third major shift was initiated by the development of third-generation sequencing technologies, like:

Pacific Biosciences (PacBio SMRT), generating reads up to 50,000 bp,

Oxford Nanopore Technologies (ONT), producing ultra-long reads exceeding 1 Mb.

While these reads have **higher error rates (5–15%)**, they enable:

Bridging large repeats, structural variants, and telomeres,

Assembly of complete bacterial genomes in a single contig,

Resolving complex loci in plant, animal, and human genomes.

Assemblers like Canu, Flye, Shasta, and wtdbg2 employ OLC-style or string graph models, which work better with long reads due to their lower reliance on fine overlap distinctions. These tools use advanced error correction techniques such as:

**Self-correction** (overlapping noisy reads with each other),

Polishing with high-quality short reads using tools like Pilon or Racon.

Moreover, these tools incorporate **repeat-aware graph simplification**, **adaptive k-mer weighting**, and **consensus generation** through multiple sequence alignments or Hidden Markov Models (HMMs).

## Stage 4: Hybrid Assemblers and Integrative Frameworks

Given the complementary strengths of short and long reads — high accuracy vs. long-range continuity — the bioinformatics community embraced **hybrid assembly** approaches. These pipelines typically:

Use long reads for layout and short reads for correction/polishing,

Leverage graph structures like hybrid De Bruijn graphs or repeat graphs,

Tools: MaSuRCA, hybridSPAdes, Unicycler, and OPERA-MS.

Such methods have demonstrated **higher contiguity (N50)** and **lower misassembly rates**, particularly in microbial and metagenomic assemblies.

#### Stage 5: Machine Learning and Beyond

Recently, **machine learning (ML)** and **deep learning** have begun influencing genome assembly. Tools like **DeepVariant** (Google) and **PEPPER-Margin-DeepVariant** (for ONT reads) improve variant calling by learning sequence patterns and systematic biases. Emerging **AI-based assemblers** aim to learn assembly rules from large datasets, adaptively handling error-prone or low-coverage regions.

Moreover, genome assembly is being transformed by:

#### Lattice deconvolution methods,

Graph-theoretic innovations, such as variation graphs and pan-genome assemblies,

Cloud-native and distributed frameworks for real-time analysis of massive datasets.

In conclusion, genome assembly tools have evolved from manual, low-throughput, overlap-based methods to highly automated, error-resilient, AI-augmented frameworks. Each generation reflects a careful trade-off between the computational complexity, sequencing platform limitations, and biological accuracy. Understanding this evolution is crucial for selecting the appropriate tools for specific genome projects — whether assembling a bacterial plasmid, a cancer genome, or an entire plant pangenome.

#### 2. Short-Read and Long-Read Assemblers:

The choice between short-read and long-read genome assemblers largely depends on the sequencing technology used and the biological complexity of the genome being studied. Both categories of assemblers are grounded in distinct computational philosophies and offer different trade-offs in terms of accuracy, completeness, contiguity, and resource requirements.

#### **Short-Read Assemblers:**

Short-read assemblers are designed to process data generated by **second-generation sequencing platforms**, such as **Illumina**, which produce high-accuracy reads typically ranging from 50 to 300 base pairs in length. These reads are economical and highly reliable in terms of base calling but are **too short** to span most repetitive elements or structural variants in complex genomes.

#### **Key Short-Read Assemblers:**

**Velvet:** One of the earliest De Bruijn graph-based assemblers. Velvet constructs a graph from k-mers and applies simplification strategies like tip clipping and bubble popping to resolve ambiguities. It is efficient for microbial genomes but struggles with eukaryotic complexity.

**SOAPdenovo:** Designed for large-scale de novo assembly of plant and animal genomes. It features a modular design with improved scaffolding algorithms and can process large genomes with high coverage. **SPAdes:** A more advanced assembler that uses multi-k De Bruijn graphs to capture information at multiple

resolution levels. It is particularly useful for single-cell and metagenomic datasets due to its built-in error correction and high assembly accuracy.

#### **Strengths of Short-Read Assemblers:**

High base accuracy: Ideal for SNP and small indel detection.

**Cost-effective sequencing:** Suitable for large-scale studies and population genomics.

Fast assembly times with relatively low computational demands for microbial genomes.

**Limitations:** 

Limited contiguity: Inability to span large repeats leads to fragmented assemblies (low N50).

Ambiguity in repeat-rich regions: Often results in unresolved scaffolds and misassemblies.

Not suitable for assembling genomes with extensive structural variations or polyploidy.

#### **Long-Read Assemblers:**

Long-read assemblers are tailored for third-generation sequencing technologies such as **Pacific Biosciences** (**PacBio**) and **Oxford Nanopore Technologies** (**ONT**). These platforms produce reads ranging from 10 kb to several megabases, offering the ability to span complex regions of the genome. However, the trade-off is **higher error rates**, especially in ONT reads.

## **Key Long-Read Assemblers:**

Canu: Successor to the Celera Assembler, designed specifically for noisy long reads. It uses a multi-stage pipeline: read correction, trimming, and assembly. It performs well with high-error data and large genomes but requires substantial computational resources.

**Flye:** Uses repeat graphs and an adaptive k-mer approach to efficiently assemble genomes with high contiguity. It supports assembly of both circular and linear genomes and is well-suited for microbial and small eukaryotic genomes.

**wtdbg2:** Also known as Redbean, this assembler converts noisy long reads into a fuzzy De Bruijn graph for efficient and ultra-fast assembly. It is optimized for speed and has been shown to produce good results on real-time ONT data.

#### **Strengths of Long-Read Assemblers:**

**High contiguity:** Capable of assembling entire chromosomes or circular genomes in single contigs.

Repeat resolution: Long reads bridge over tandem and interspersed repeats, improving genome completeness.

**Structural variant detection:** Ideal for identifying large insertions, deletions, inversions, and duplications. **Limitations:** 

**Higher error rates:** Raw long reads (especially from ONT) require extensive polishing using tools like **Pilon, Medaka,** or **Racon**.

**Expensive and resource-intensive:** Long-read sequencing is costlier and requires higher computational memory and CPU time.

**Coverage bias:** May produce uneven coverage in GC-rich or AT-rich regions, affecting assembly fidelity. **Contextual Use Cases:** 

**Microbial Genomics:** Short-read assemblers like SPAdes are often sufficient for bacteria, while long-read tools like Flye provide complete circular assemblies in a single contig.

**Eukaryotic Genomes:** Complex genomes with large repeats and heterozygosity demand long-read assemblers (e.g., Canu or wtdbg2) for accurate chromosome-scale assemblies.

**Cancer Genomics:** Long-read assemblers are preferred for resolving tumor heterogeneity and detecting structural variants that short-read assemblers might miss.

**Metagenomics:** SPAdes (metaSPAdes) works well for high-complexity microbial communities, but hybrid approaches (e.g., using Flye for dominant species) may offer better resolution.

Short-read and long-read assemblers each have distinct strengths and weaknesses, and the optimal choice depends on project goals, genome complexity, and available resources. While short-read assemblers offer high accuracy and affordability for simpler genomes, long-read assemblers unlock the full potential of modern genomics by delivering more complete, contiguous, and structurally accurate assemblies. Increasingly, **hybrid strategies** that combine the strengths of both technologies are being adopted to achieve balanced, high-quality results.

#### 3. Hybrid Genome Assembly Approaches:

Hybrid genome assembly represents a powerful strategy that synergizes the **high accuracy of short reads** (e.g., from Illumina) with the **long-range continuity of long reads** (e.g., from Oxford Nanopore or PacBio). This approach has emerged as a practical solution to overcome the limitations inherent in using either technology alone. While short reads excel at base-level accuracy, they fall short in resolving repeats and large-scale structural features. Conversely, long reads bridge such repetitive or complex regions but often introduce base errors. Hybrid assemblers integrate both data types to produce **highly contiguous**, **complete, and accurate genome assemblies**, particularly in complex eukaryotic or polyploid genomes.

## Why Combine Illumina with Nanopore/PacBio Reads?

Illumina short reads are known for their high accuracy (error rate < 0.1%) and low cost, making them excellent for detecting single-nucleotide polymorphisms (SNPs) and small insertions or deletions.

Nanopore and PacBio long reads, despite having higher error rates (initially 5–15%, now reduced via consensus polishing), can span repetitive elements, structural variants, and GC-rich or AT-rich regions, which are problematic for short-read assemblers.

By integrating both:

The **long reads define the scaffold or backbone** of the genome.

The **short reads are used to polish or correct base-level errors**, significantly enhancing the accuracy of the final assembly.

#### **Key Hybrid Assembly Tools:**

MaSuRCA (Maryland Super Read Cabog Assembler):

Combines the efficiency of short-read De Bruijn graphs with the resolving power of long-read scaffolding.

It creates **super-reads** by extending short reads into longer contigs using overlapping k-mers.

Long reads are then used to link, scaffold, and resolve repeats, with final consensus derived from the combined read data.

Particularly effective for large eukaryotic genomes and has been used in **plant and vertebrate assemblies**. **hybridSPAdes:** 

An extension of the **SPAdes assembler**, specifically designed to incorporate both **Illumina short reads** and long reads (ONT or PacBio).

Constructs an initial graph using short reads, then uses long reads to **extend and resolve ambiguous paths** in the graph.

Well-suited for small bacterial genomes and single-cell studies.

Includes robust error correction modules and multiple k-mer strategies for increased resolution.

#### Unicycler:

Designed primarily for bacterial genome assembly, Unicycler integrates short-read accuracy with long-read completeness.

Starts with a short-read assembly (e.g., using SPAdes), then **bridges contigs** using long reads to form **complete circular genomes**.

It's one of the few tools that can produce single-contig assemblies of bacterial plasmids and chromosomes.

Automatically performs **read polishing** (e.g., using Racon or Pilon), resulting in a final assembly with **low error rates and high structural integrity**.

#### **Benefits of Hybrid Assembly:**

#### **Improved Contiguity and Scaffold Length:**

Long reads provide the necessary coverage to connect distant contigs and resolve gaps that short reads cannot span.

Assemblies often show a higher N50 and lower number of contigs, indicating greater completeness.

## **Resolution of Repeats and Structural Variants:**

Repetitive regions, such as **transposable elements**, **telomeres**, and **centromeres**, can be properly assembled because long reads span them.

Enables accurate detection of large insertions, deletions, duplications, and inversions.

#### **Enhanced Base-Level Accuracy:**

Polishing with short reads corrects the base-calling errors typically present in raw long reads.

Tools like Pilon or Medaka are often used after hybrid assembly for fine-tuning.

#### **Cost-Effectiveness:**

Using a modest number of long reads in combination with cheap short-read data provides a budget-friendly path to high-quality assemblies, especially for small to medium-sized genomes.

#### **Limitations and Considerations:**

**Computational Demand:** Hybrid assembly requires large memory and CPU resources, especially during graph construction and polishing stages.

**Read Normalization and Quality Control:** Preprocessing of both short and long reads (adapter removal, quality filtering, error correction) is crucial to prevent misassemblies.

**Tool Compatibility:** Not all hybrid assemblers handle every data format uniformly; the choice of assembler may depend on the specific read technologies and genome type.

Hybrid genome assembly has become an essential strategy in modern genomics, enabling researchers to capitalize on the complementary strengths of Illumina and Nanopore/PacBio sequencing technologies. Tools like MaSuRCA, hybridSPAdes, and Unicycler exemplify how integration of data types can produce assemblies that are not only structurally accurate but also highly polished. As sequencing costs decline and algorithmic innovations continue, hybrid assembly is likely to remain a mainstay in microbial genomics, clinical diagnostics, evolutionary biology, and complex genome projects, providing the best of both sequencing worlds.

#### 4. Error Correction and Assembly Polishing:

Genome assembly is rarely perfect in its first pass, especially when using raw reads from long-read technologies such as **Oxford Nanopore** and **PacBio**, which, despite offering long-range continuity, are prone to **base-calling errors**, **insertions/deletions (indels)**, and **structural inaccuracies**. To achieve a high-quality and biologically interpretable genome—especially one suitable for **clinical applications or diagnostic use**—a critical post-processing step called **assembly polishing** is employed. This step enhances **base-level accuracy**, improves **structural fidelity**, and corrects **algorithmic artifacts** left over from the initial assembly.

## What is Error Correction and Assembly Polishing?

**Error correction** typically happens at the read level, before assembly. It improves the quality of the raw sequencing data by removing or modifying erroneous reads.

**Assembly polishing**, on the other hand, happens after an initial draft genome is constructed. It aligns the original reads (short or long) back to the assembly and **detects mismatches**, **gaps**, **or misjoins**.

This process ensures that the final assembly **faithfully reflects the true biological sequence** and meets the high accuracy requirements of downstream applications such as **gene annotation**, **variant detection**, and **comparative genomics**.

#### **Key Tools for Assembly Polishing:**

#### **Pilon (Broad Institute):**

Designed for polishing assemblies using accurate short reads (e.g., Illumina).

Works by mapping Illumina reads to the draft genome using aligners like BWA or Bowtie2.

Detects:

**Base substitutions** (A $\rightarrow$ G, T $\rightarrow$ C, etc.)

Small insertions and deletions (indels)

Misassembled regions, including collapsed repeats or wrongly joined contigs

Strengths: Produces extremely accurate final sequences, suitable for bacterial genomes, fungal genomes, and even small eukaryotes.

Use case: Pilon is often used in hybrid assemblies where long-read contigs are corrected with short-read alignments.

#### Racon:

Designed for use with long-read technologies, such as ONT or PacBio.

Polishes assemblies using the same long reads that generated the assembly.

Employs a partial order alignment (POA) algorithm to realign reads and create a corrected consensus.

Racon is often run in **multiple rounds** for better results and is typically paired with tools like **Minimap2** for fast read mapping.

**Strengths**: Fast, memory-efficient, and highly scalable, making it ideal for polishing large genomes using long reads alone.

#### Medaka (Oxford Nanopore Technologies):

A neural network-based polisher specifically optimized for **Oxford Nanopore** reads.

After an initial Racon polish, Medaka uses **machine learning models** to further refine base calls, especially around **homopolymeric regions** (e.g., AAAAA... or TTTTT...), which are known to be problematic for ONT base-callers.

Unlike traditional tools, Medaka does not rely on standard alignment methods alone—it interprets sequencing signal features, improving **read-context-based correction**.

**Strengths**: Outperforms traditional tools for **ONT data**, providing polishing quality comparable to Illumina-based tools, especially for **clinical-grade human genome assemblies**.

#### What Types of Errors Are Corrected?

**Base Substitution Errors**: Incorrect single nucleotide identification (e.g., C instead of T).

**Insertions/Deletions (Indels)**: Especially problematic in repetitive sequences and homopolymers; can cause frameshifts in gene coding regions.

**Misassemblies**: Incorrect joining of unrelated contigs or collapse of repetitive regions; polishing helps realign and restructure such regions.

**Consensus Errors**: Mistakes in consensus base calling during initial assembly that are corrected through iterative polishing rounds.

#### Importance of Polishing for Clinical-Grade Assemblies:

Genome assemblies intended for diagnostic or therapeutic applications, such as detecting disease-causing mutations, pharmacogenomics profiling, or pathogen surveillance, require:

>99.9% accuracy (often referred to as Q40 or higher)

Accurate representation of regulatory and coding regions

No false-positive variant calls, which can mislead clinical decisions

Without proper polishing, even a few sequencing or assembly errors can lead to:

#### Misannotation of genes

#### Erroneous identification of SNPs or structural variants

#### Inaccurate phylogenetic analyses or population studies

Therefore, tools like **Pilon**, **Racon**, and **Medaka** are indispensable in transforming a draft assembly into a **high-fidelity reference** genome.

#### **Typical Polishing Pipelines:**

For example, a hybrid polishing strategy for bacterial genome assembly might look like:

Assemble long-read data using Flye or Canu

Run two rounds of Racon using long reads

Apply Medaka for deep ONT-based correction

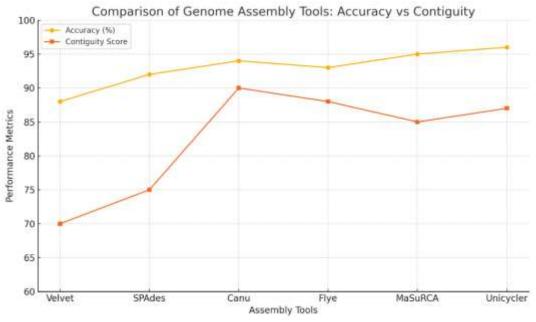
Polish with Pilon using high-coverage Illumina reads for base-level refinement

This process yields an assembly that is both **structurally contiguous** and **base-wise accurate**, ready for publication or clinical interpretation.

Error correction and assembly polishing are vital stages in any genome assembly pipeline, especially in projects requiring high-confidence sequences. Tools like **Pilon**, **Racon**, and **Medaka** each offer unique strengths tailored to specific sequencing platforms. Their combined use not only improves **base accuracy** and **genome integrity**, but also ensures the final assemblies meet the **stringent standards of scientific rigor and clinical utility**. As sequencing technologies continue to evolve, so too will polishing tools, integrating **deep learning** and **real-time feedback** mechanisms for even greater precision.

# Comparison of Genome Assembly Tools: Accuracy vs

## Contiguity



#### **Summary:**

Advancements in bioinformatics tools for genome assembly have significantly enhanced our ability to decode complex genomes with higher accuracy and efficiency. From early graph-based algorithms to modern hybrid and AI-powered tools, the field has grown to accommodate a range of sequencing technologies and genome complexities. Hybrid assembly methods bridge the gap between accuracy and read length, while error correction tools ensure reliability for downstream applications. Looking forward, machine learning and scalable computing infrastructure will play vital roles in meeting the increasing demands of genome science. The continuous evolution of these tools will not only support research in basic and applied genomics but also accelerate breakthroughs in personalized medicine and biotechnology.

#### **References:**

- Bankevich, A., et al. (2012). SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing. Journal of Computational Biology, 19(5), 455–477.
- Zimin, A. V., et al. (2013). The MaSuRCA genome assembler. Bioinformatics, 29(21), 2669–2677.
- Li, H. (2016). Minimap and miniasm: fast mapping and de novo assembly for noisy long sequences. Bioinformatics, 32(14), 2103–2110.
- Vaser, R., et al. (2017). Fast and accurate de novo genome assembly from long uncorrected reads. Genome Research, 27(5), 737–746.
- Walker, B. J., et al. (2014). Pilon: An Integrated Tool for Comprehensive Microbial Variant Detection and Genome Assembly Improvement. PLOS ONE, 9(11), e112963.
- Wick, R. R., et al. (2017). Unicycler: Resolving bacterial genome assemblies from short and long sequencing reads. PLOS Computational Biology, 13(6), e1005595.
- Kolmogorov, M., et al. (2019). Assembly of long, error-prone reads using repeat graphs. Nature Biotechnology, 37(5), 540–546.
- Sahlin, K., & Medvedev, P. (2021). Error correction enables use of Oxford Nanopore technology for reference-free transcriptome analysis. Nature Communications, 12, 2–12.
- Poplin, R., et al. (2018). A universal SNP and small-indel variant caller using deep neural networks. Nature Biotechnology, 36(10), 983–987.
- Wenger, A. M., et al. (2019). Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. Nature Biotechnology, 37(10), 1155–1162.
- Shafin, K., et al. (2021). Haplotype-aware variant calling with PEPPER-Margin-DeepVariant enables high accuracy in nanopore sequencing. Nature Methods, 18(11), 1322–1332.
- Logsdon, G. A., Vollger, M. R., & Eichler, E. E. (2020). Long-read human genome sequencing and its applications. Nature Reviews Genetics, 21(10), 597–614.