# Genome assembly strategies
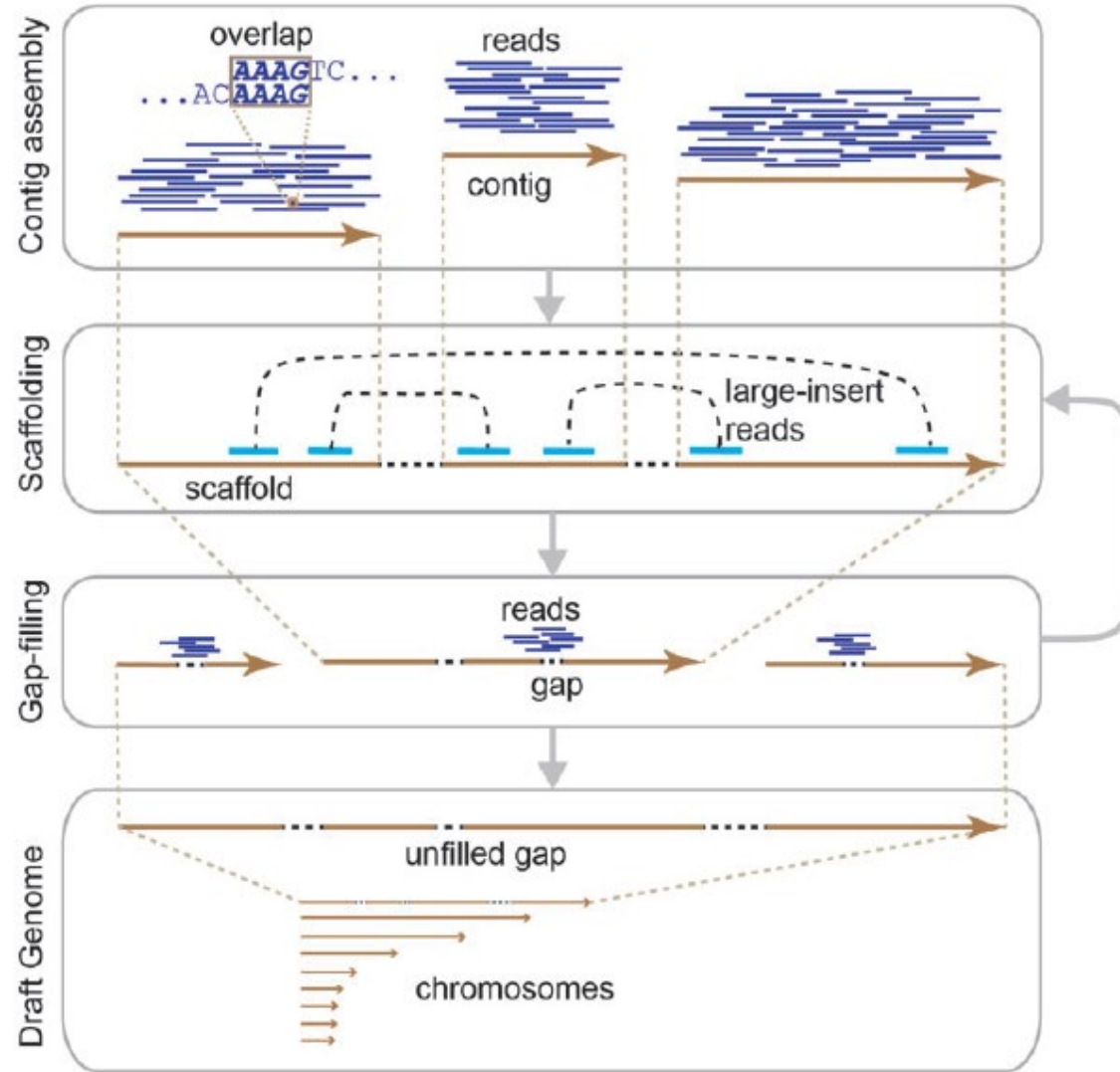
Arturo Vera Ponce de Leon

May 2019

veraponcedeleon.1@osu.edu
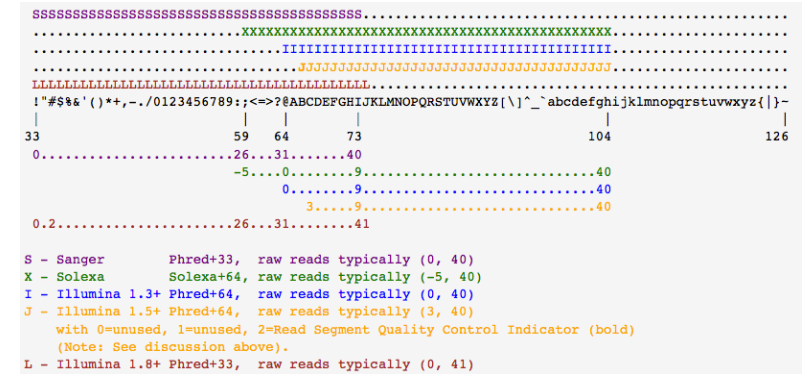
# Genome assembly

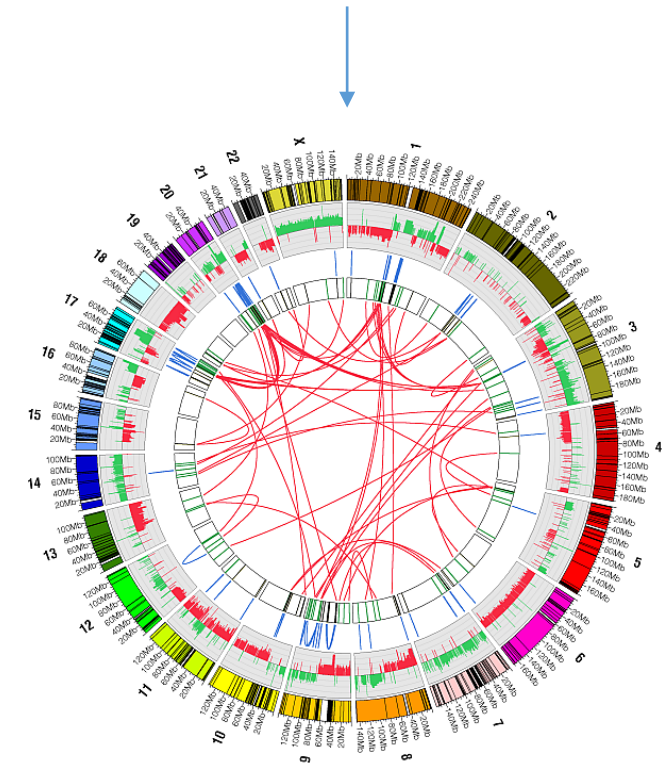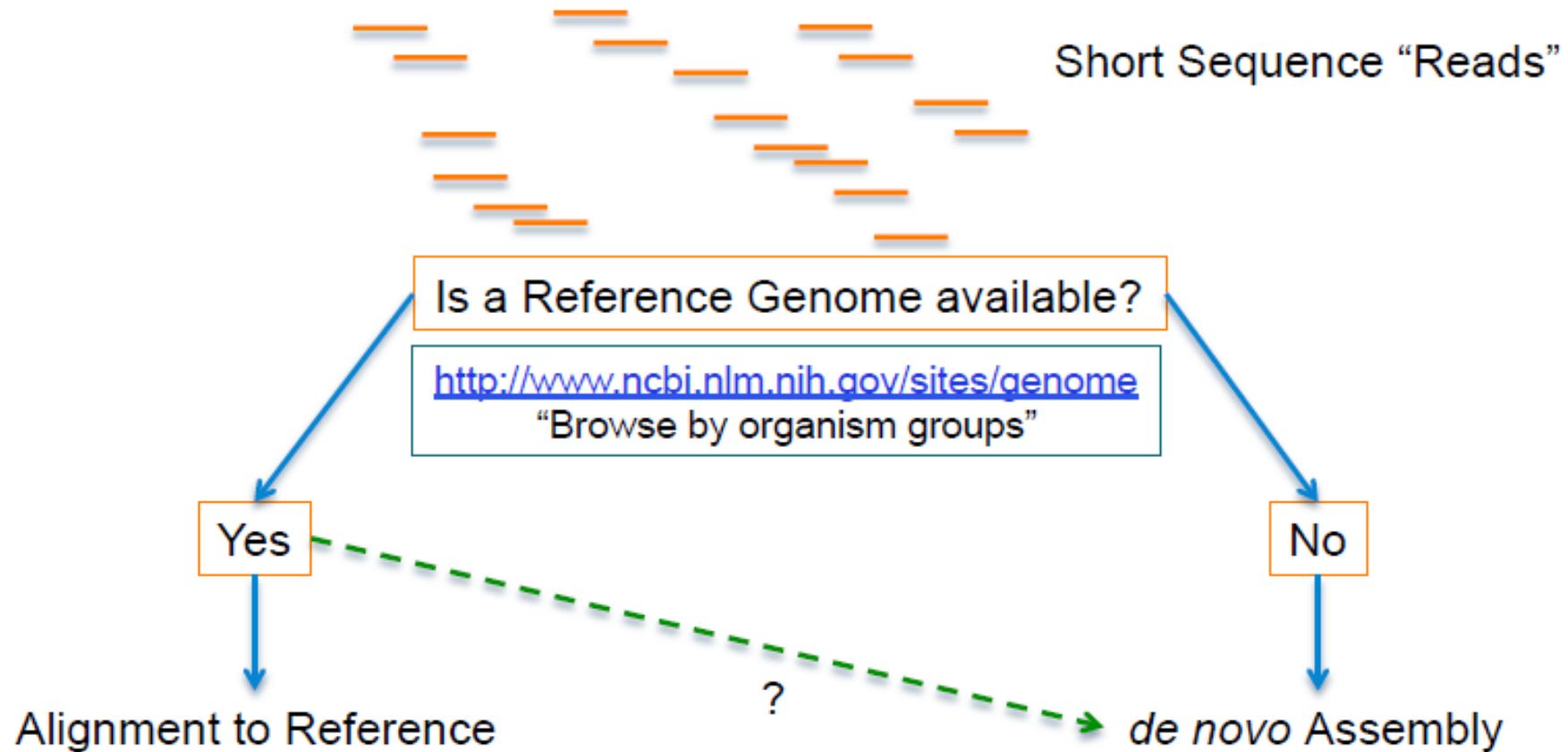# Genome assembly



Reads
(fastq)

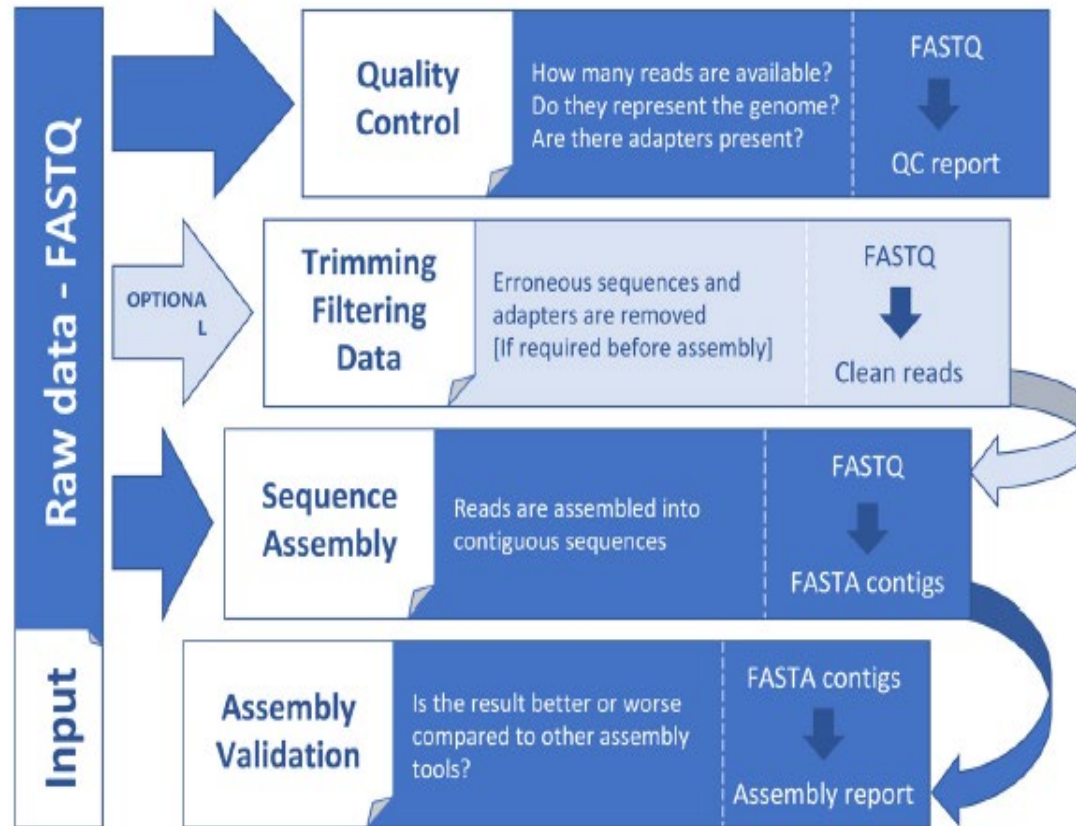Summarizing

Chromosome

# Genome *de-novo* assembly or reference mapping

# Genome Assembly

## Work flow to classic genome assembly strategies

## Software used in this lecture



Figure 2. General steps in a genome assembly workflow. Input and output data are indicated for each step.

**fastQC**

**TrimGalore**

**IDBA**
**SPADES**

**QUAST**
**BUSCO**
**CheckM**

# Strategies to genome assembly

- Algorithms

1. Greedy

2. Overlap-layout-
   consensus (OLC)

3. De Bruijn Graph



**A** Read Layout

```
R₁:  GACCTACA
R₂:    ACCTACAA
R₃:     CCTACAAG
R₄:      CTACAAGT
A:        TACAAGTT
B:         ACAAGTTA
C:          CAAGTTAG
X:         TACAAGTC
Y:          ACAAGTCC
Z:           CAAGTCCG
```

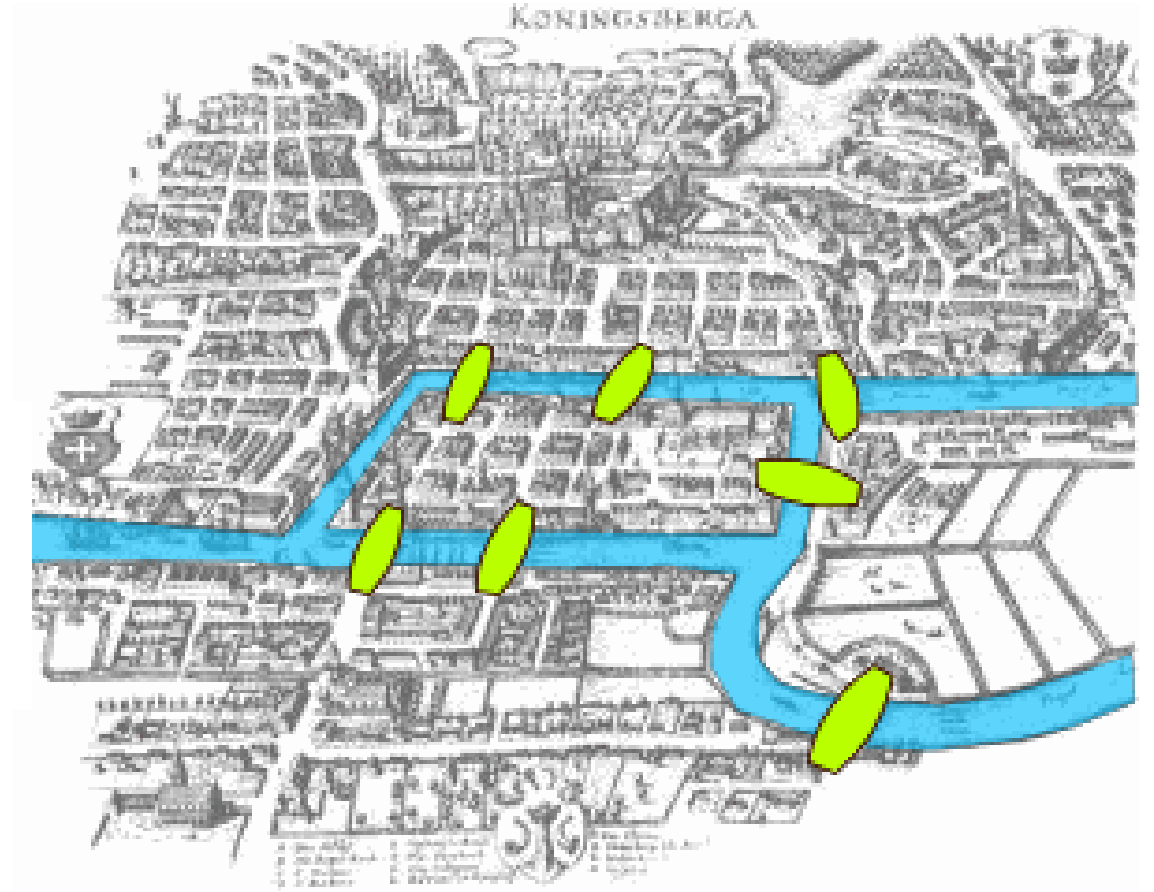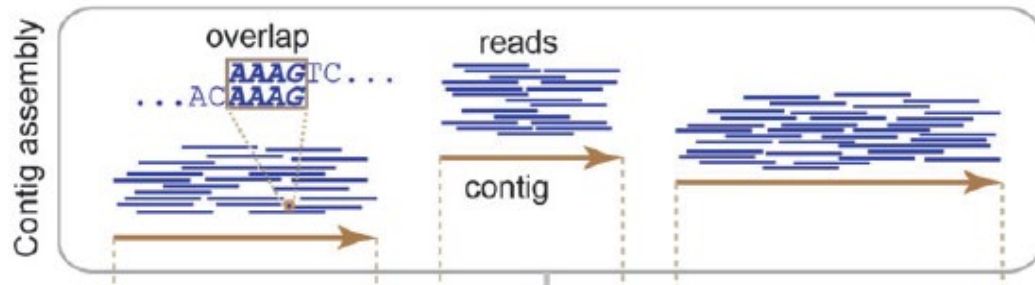**B** Overlap Graph

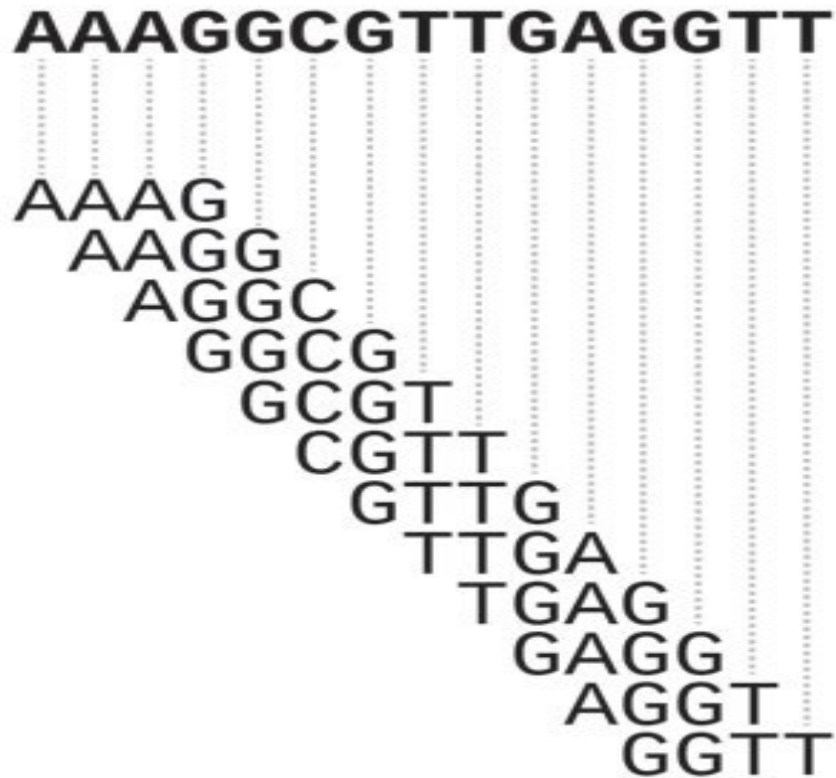**C** de Bruijn Graph

# Steeps to genome assembly using De-Bruijn graph

**The basic strategy for de novo assembly for short NGS reads comprises three steps: (i) contig assembly, (ii) scaffolding and (iii) gap filling.**





**Seven Bridges of Königsberg**

# The K-mers: divide and conquer
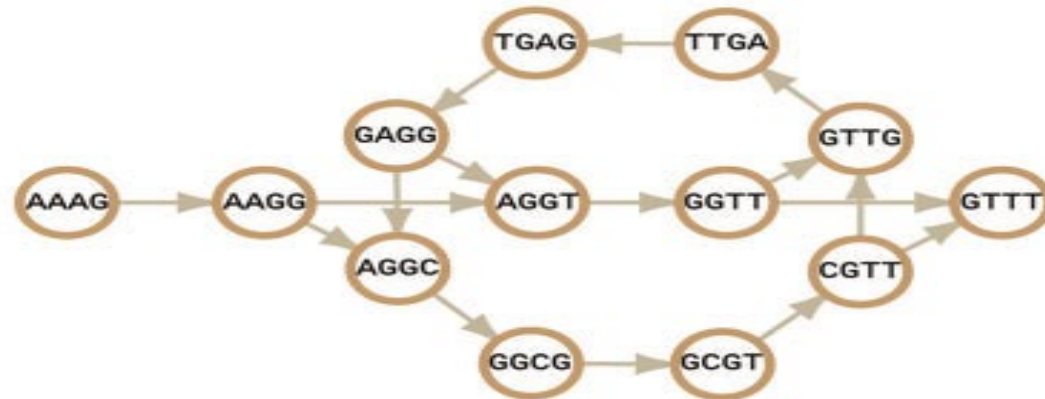


**A** Short read to *k*-mers (*k*=4)

**AAAGGCGTTGAGGTT**

AAAG
AAGG
 AGGC
  GGCG
   GCGT
    CGTT
     GTTG
      TTGA
       TGAG
        GAGG
         AGGT
          GGTT

**B** Eulerian de Bruijn graph

**C** Hamiltonian de Bruijn graph

# The K-mers: divide and conquer

• It breaks reads into <span style="color:red">successive</span> k-mers and the graph maps the k-mers

• Each k-mer is a node and edges are drawn between each k-mer in a read.

• Repeat sequences create a fork in the graph; alternative sequences create a bubble.

• The k-mer size can only be determined by "trial and error".

• A small value of K will create a complex graph but a large value of K may miss small overlaps. A good starting point would be a k-mer size that is 2/3 the size of the read

• Good for short reads or small genomes. With long reads and/or large genomes, may require lots of RAM (e.g., ~0.5 TB for human)

# Let´s go to assembly some bacterial genomes

**IDBA-UD: a *de novo* assembler for single-cell and metagenomic sequencing data with highly uneven depth**

Yu Peng, Henry C. M. Leung*, S. M. Yiu and Francis Y. L. Chin
Department of Computer Science, The University of Hong Kong, Pokfulam Road, Hong Kong
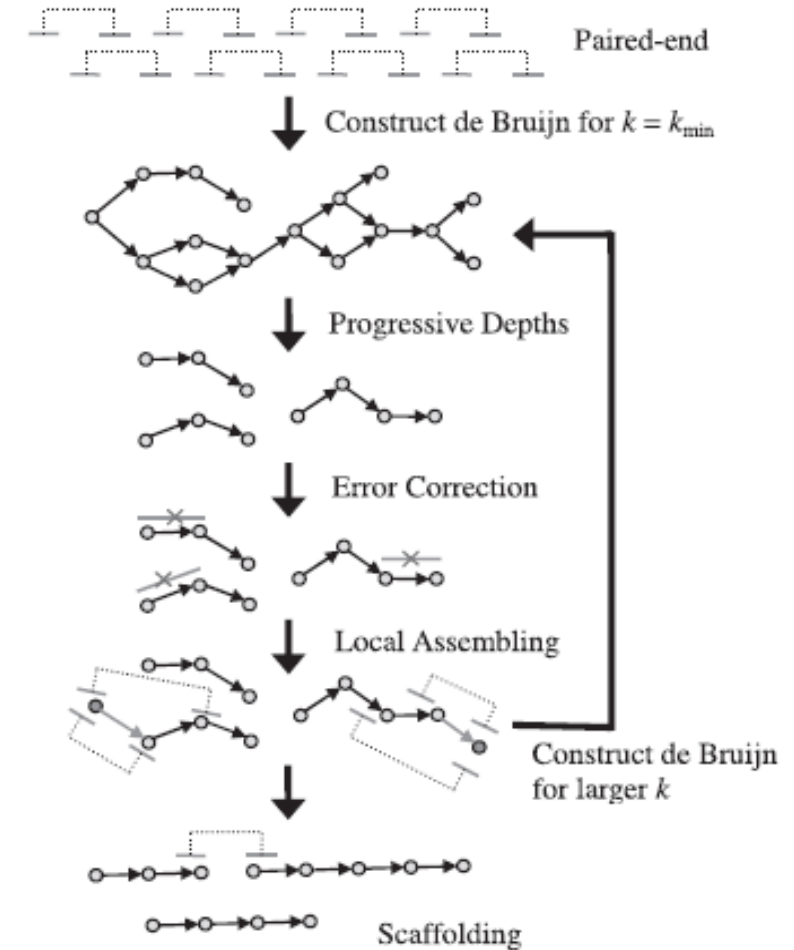Associate Editor: Michael Brudno

**Fig. 1.** Flowchart of IDBA-UD

# Evaluating the assembly

§ **Genome assembly results:**
• **contig size and number of contigs produced**
• **scaffold size and number**
• **N50 and N90**
§ **Coverage**
§ **GC Content**
§ **Genome annotation**
• repeats analysis and annotation
• protein-coding gene annotation (including gene structure prediction and gene function annotation)
• non-coding RNA gene annotation (including annotation of microRNA, tRNA, rRNA, and other ncRNA)
• transposon and tandem repeats annotation
§ Comparative genomics and evolution (chromosome structure, conserved gene families)

# Basic stats

Basic statistics
**N50** the length of the shortest contig such that the sum of contigs of equal length or longer is at least 50% of the total length of all contigs.

Contig size (bp)
3000
2000 N50
1200
800
600 N90
400
Total: **8000**

N90 = the length of the shortest contig such that the sum of contigs of equal length or longer is at least 90% of the total length of all contigs.

# SPADEs

Original Articles

## SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing

ANTON BANKEVICH,[1,2] SERGEY NURK,[1,2] DMITRY ANTIPOV,[1] ALEXEY A. GUREVICH,[1] MIKHAIL DVORKIN,[1] ALEXANDER S. KULIKOV,[1,3] VALERY M. LESIN,[1] SERGEY I. NIKOLENKO,[1,3] SON PHAM,[4] ANDREY D. PRJIBELSKI,[1] ALEXEY V. PYSHKIN,[1] ALEXANDER V. SIROTKIN,[1] NIKOLAY VYAHHI,[1] GLENN TESLER,[5] MAX A. ALEKSEYEV,[1,6] and PAVEL A. PEVZNER[1,4]