# Next-generation sequencing

Lecture 6

# Assembly Algorithms

Main algorithm used:
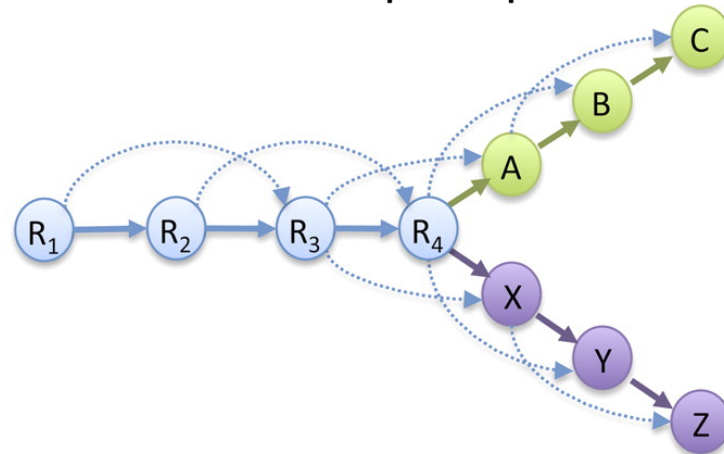
- Greedy algorithms
- Overlap Layout Consensus
- De bruijn graphs

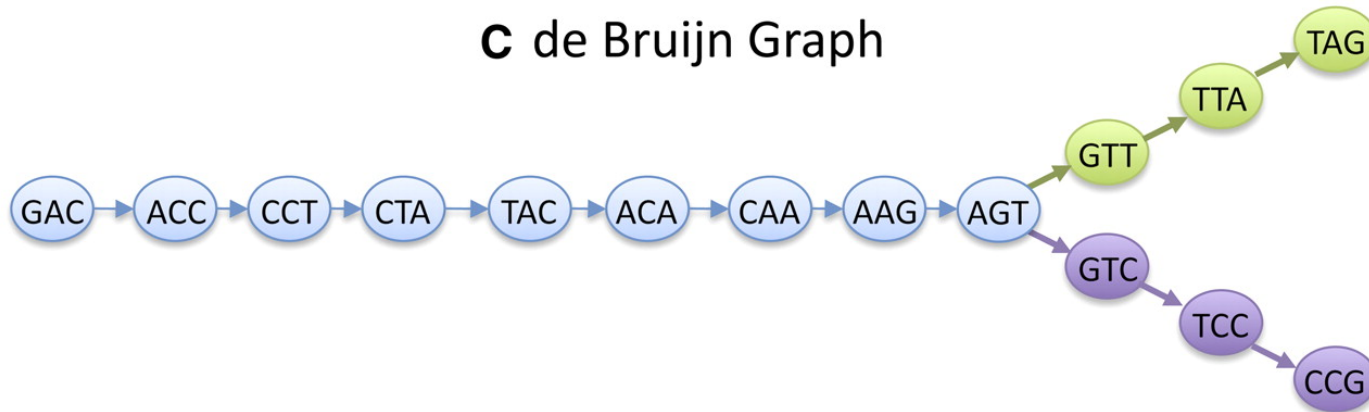# Differences between an overlap graph and a de Bruijn graph for assembly.

## A Read Layout

R$_1$: GACCTACA
R$_2$:   ACCTACAA
R$_3$:     CCTACAAG
R$_4$:       CTACAAGT
A:          TACAAGTT
B:            ACAAGTTA
C:              CAAGTTAG
X:          TACAAGTC
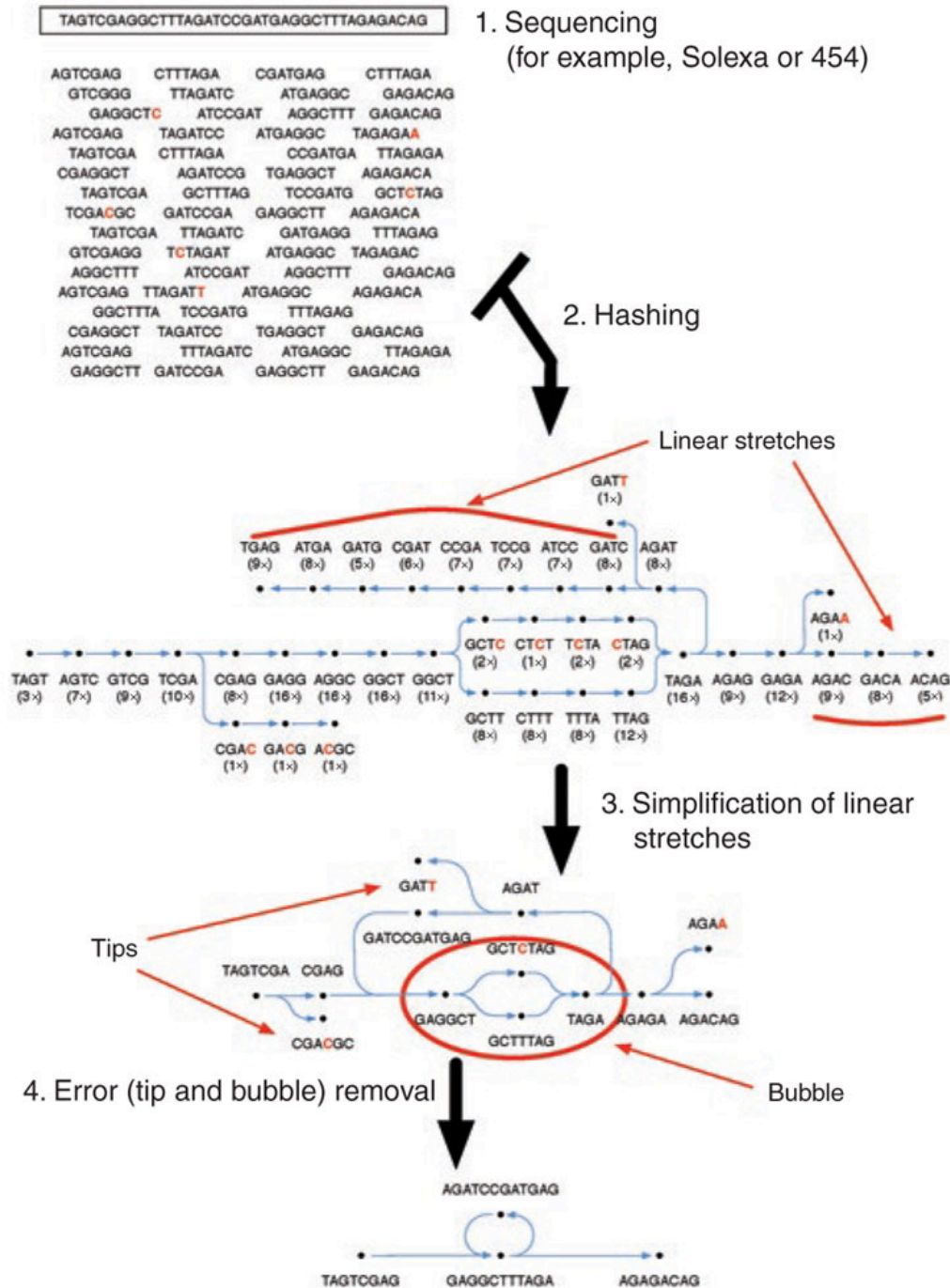Y:            ACAAGTCC
Z:              CAAGTCCG

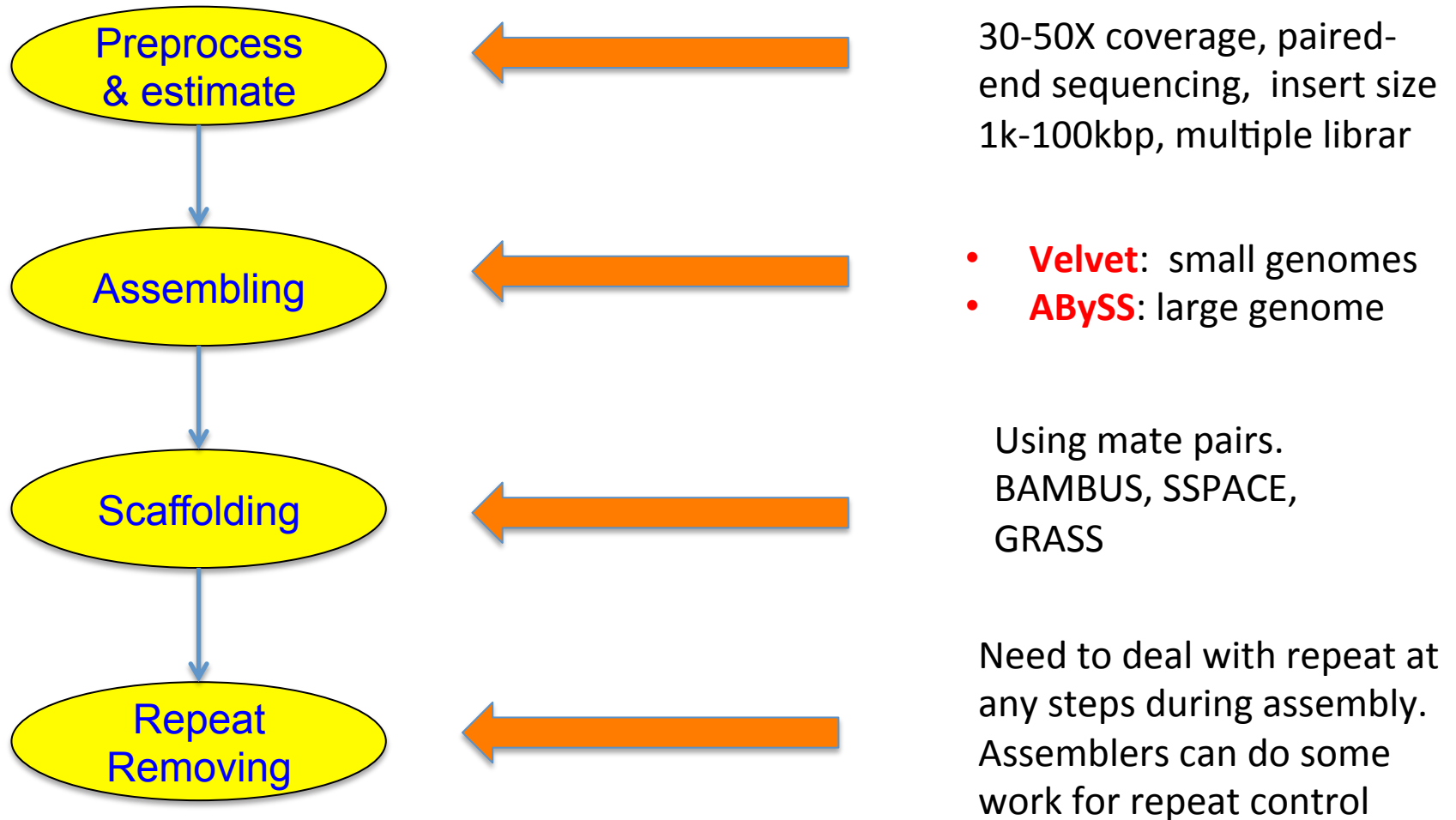## B Overlap Graph

## C de Bruijn Graph

# De Bruijn Graphs

1. Get k-bp (k-mer) subsequences for reads.
2. k-mers in the reads are collected into nodes and the coverage at each node is recorded. Link two k-mer nodes if they have overlap.
3. the graph is simplified to combine nodes that are associated with the continuous linear stretches into single, larger nodes of various k-mer sizes.
4. error correction removes the tips and bubbles that result from sequencing errors and creates a final graph structure that accurately and completely describes in the original genome sequence.

Flicek, Nature Methods, 2009

# De Bruijn Assemblers

- Euler: http://nbcr.sdsc.edu/euler/ , Sanger, 454, 2001-2006

- **Velvet**: http://www.ebi.ac.uk/~zerbino/velvet/, small genomes, Sanger, 454, Solexa, SOLiD, 2007-2009 (very good for small genome)

- **ABySS**: http://www.bcgsc.ca/platform/bioinfo/software/abyss, large genome, Solexa, SOLiD, 2008-2011 (for very large genome)

- **SOAP-denovo**: http://soap.genomics.org.cn/soapdenovo.html, Solexa, 2009

- ALLPATH-LG: http://www.broadinstitute.org/software/allpaths-lg/blog/, large genome, Solexa, SOLiD, 2011 (very good performance bu require 2 lib of different insert sizes)

- IDBA-UD: http://i.cs.hku.hk/~alse/hkubrg/projects/idba_ud/, Sanger, 454,Solexa, 2010 (metagenomic, doesn't rely on coverage to remove error)

# Assembly Pipeline



**Preprocess & estimate** — 30-50X coverage, paired-end sequencing, insert size 1k-100kbp, multiple librar

**Assembling**
- **Velvet**: small genomes
- **ABySS**: large genome

**Scaffolding** — Using mate pairs. BAMBUS, SSPACE, GRASS

**Repeat Removing** — Need to deal with repeat at any steps during assembly. Assemblers can do some work for repeat control

# Assessing Assembly Quality

- Why do we need QC?
  - Misassembly correction is expensive
  - some assemblers have a simple quality-control method that does not capture larger errors

- Common measures of quality:
  - number and sizes of contigs (N50)
    - Assumption: few large contigs is better than many small contigs.
    - True because there are less gaps in the former, but, does not account for the possibility of misassemblies.
  - And more ..
  - Compare with a complete sequence

# Assembly validation

**N50** is the most commonly used metric:

Weighted median such as 50% of your assembly is contained in contigs with length >=N50

1. Make a list L of positive integers (contig lengths).
2. Create another list L', which is identical to L, except that every element n in L has been replaced with n copies of itself.
3. The median of L' is the N50 of L.

# Assembly validation

For example:

L = {2, 2, 2, 3, 3, 4, 8, 8},

L'={2,2, 2,2, 2,2, 3,3,3, 3,3,3 ,4,4,4,4, 8,8,8,8,8,8,8,8, 8,8,8,8,8,8,8,8}

N50 of L is the median of L'.

N50=(4+8)/2 = 6.

# Assembly validation

While the N50 value thus quantifies the ability of the assembly algorithm to combine reads into large seamless blocks, it fails to capture all aspects of assembly quality.

For example, artificially high N50 values can be obtained by lowering thresholds for amalgamating smaller blocks of contiguous reads, resulting in misassembled contigs.

N50 values fail to reflect fine-scale inaccuracies, such as substitution and indel errors.
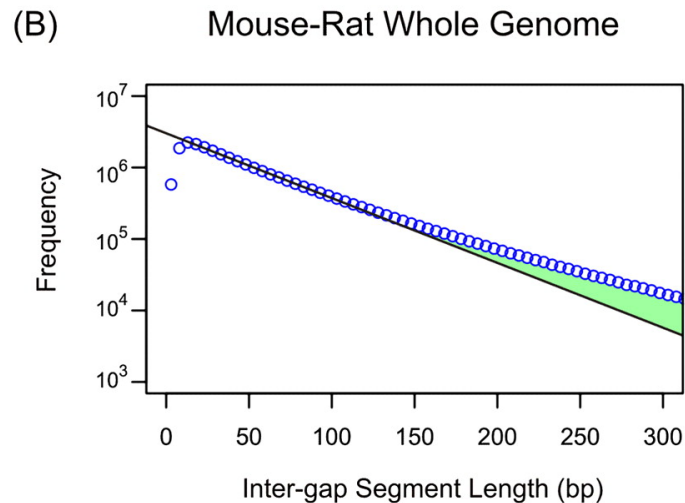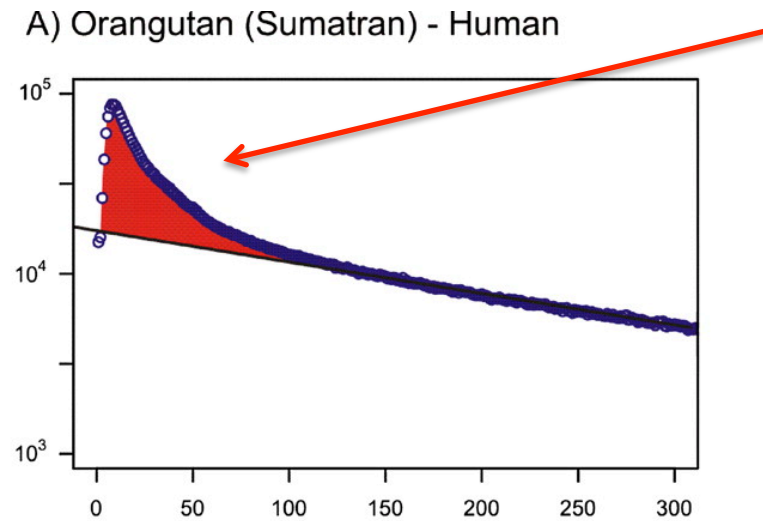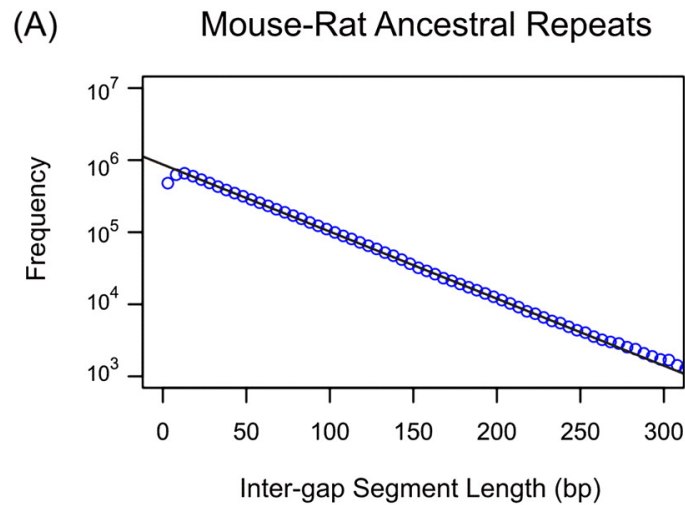
# Assembly validation

- Quality at the nucleotide level for contigs can be used to detect fine-scale inaccuracies, such as substitution and indel errors.

- Method 1: Once assembled, a base is assigned a consensus quality score (CQS) depending on its read depth and the quality of each base contributing to that position. (Huang and Madan 1999, Genome Research, 9: 868–877).

- Method 2: A multiple sequence alignment of reads is constructed and a consensus sequence along with a quality value for each base is computed for each contig.

# Assembly validation

- Method 3: a statistical and comparative genomics method that quantifies the fine-scale quality of a genome assembly and that has the merit of being complementary to the aforementioned approaches.
- This approach estimates the abundance of indel errors between aligned genome pairs, by separating these from true evolutionary indels.
- indel mutations leave a precise and determinable fingerprint on the distribution of ungapped alignment block lengths. These block lengths, which represent distances between successive indel mutations are intergap segment (IGS) lengths.

# Assembly validation



errors

(A) Mouse-Rat Ancestral Repeats

(B) Mouse-Rat Whole Genome

A) Orangutan (Sumatran) - Human

Under the **neutral indel model**, these inter-gap segment (IGS) lengths are expected to follow a geometric frequency distribution.

Meader et al., Genome Research, 2010, 20(5):675

# Assembly validation

Compare with existing genes.

CEGMA: Core Eukaryotic Genes Mapping Approach

- Looks in your assembly for genes that should be there
- Usually best assembly have best CEGMA score
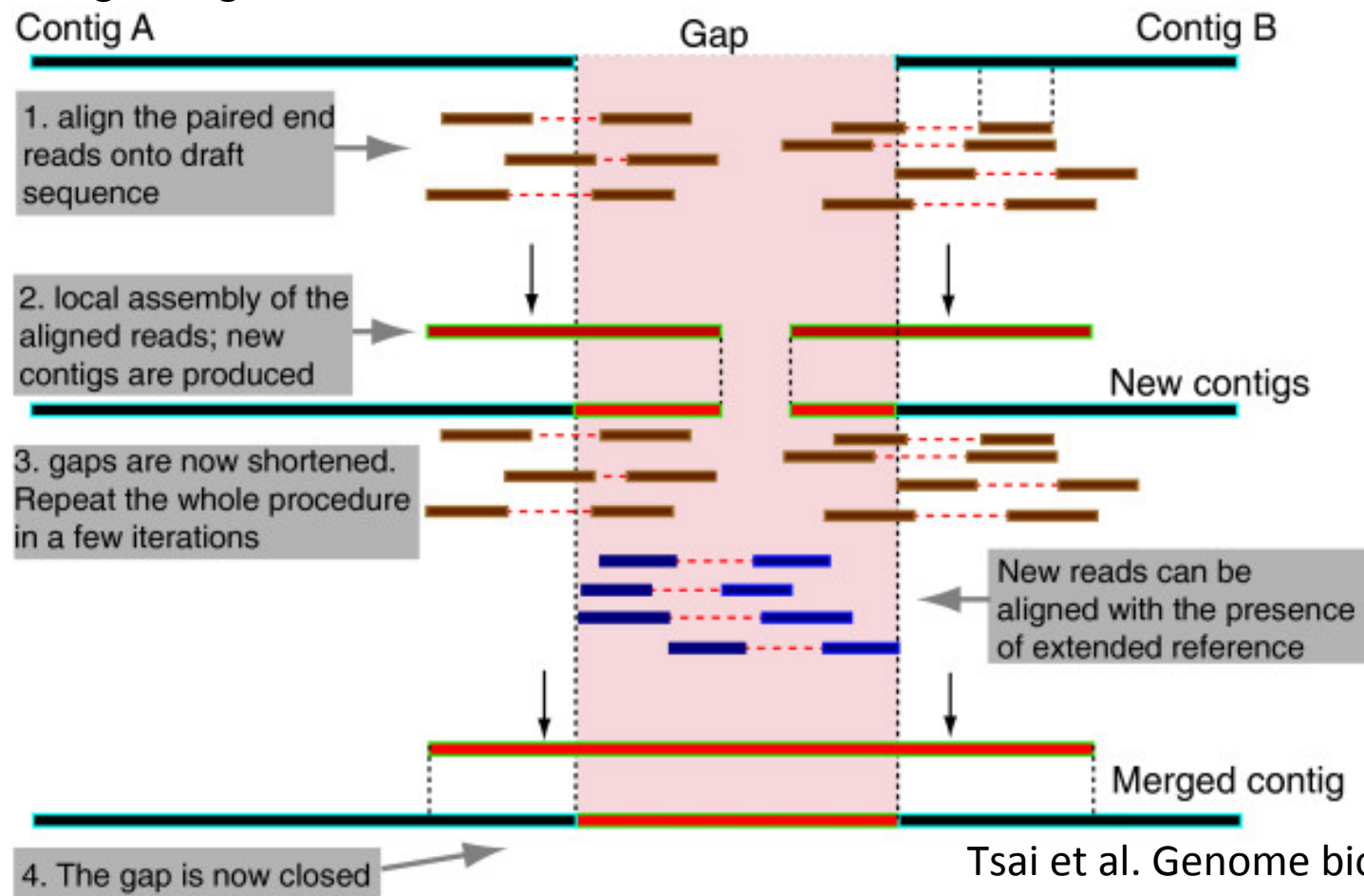
http://korflab.ucdavis.edu/datasets/cegma/

# What makes an assembly good?

- High coverage: 50 to 100X

- Different but precise insert size libraries (Paired end from different library sizes will allow you to stitch across several repeat type.)

- Avoid large number of variant.

- Sequencing errors will increase the size of the graph before correction and will sometime create branches that look real

- Error Correction: Correct the read before assembly

# What makes your assembly better?

IMAGE: Gap Filling. improve draft genome assemblies by aligning sequences against contig ends and performing local assemblies to produce gap-spanning contigs.



Tsai et al. Genome biology 2010

# Whole genome sequencing

- *De Novo* sequencing

- <span style="color:red">Mapping assembly (Reference-guided assembly) (Resequencing)</span>

  "DNA resequencing is the task of sequencing a DNA region for an individual given that a reference sequence for this region is already available for the specific species. "

  - ✓ Whole genome assembly
  - ✓ Variant discovery
    Discover or quantitate rare sequence variants
    HIV mutants within a single patient
    Scan for mutations in tumor samples