# Genome assembly post-processing
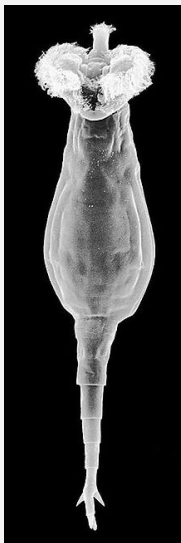
**Nadège Guiglielmoni**

# Read pre-processing

- **Adapter trimming**

- **Read filtering**: select longest/highest quality: Filtlong

- **Read correction**: reduce error rate of long reads

  - **self correction:** long reads only

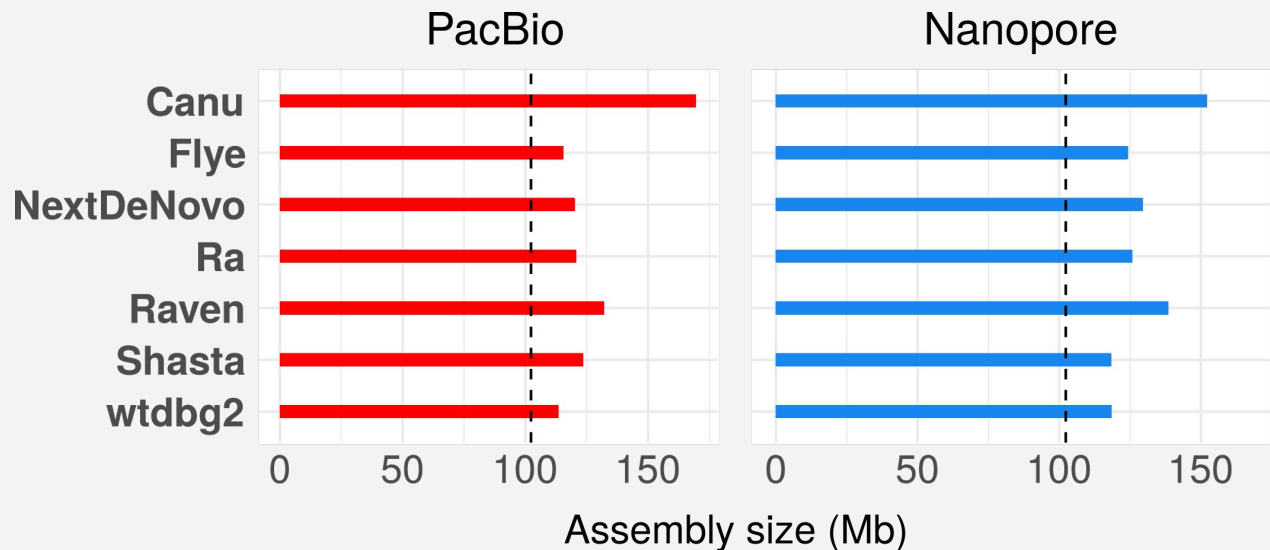  - **hybrid correction:** long reads & short reads

# Assembly post-processing

▸ **Polishing**: reduce errors

▸ **Haplotig purging**: remove uncollapsed haplotypes

▸ **Scaffolding**: increase contiguity

▸ **Gap filling**: find missing sequences

# Haplotig purging

*Adineta vaga*



Who Needs Sex (or Males) Anyway?
Liza Gross, PloS Biology, 2007

Expected haploid size 102 Mb



Assembly size (Mb)

# Haplotig purging

Haplotype 1

A T T A C C A G T C T C A A **T G G A T G G C T A C T C** T T T G A C G A T A G C T

Haplotype 2

A T T A C C A G T C T C A A **A G G C T G C T A G T G** T T T G A C G A T A G C T

**Assembly process**

**Assembly output**

Good haploid assemblies

contig 1 ✓
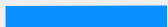
OR

contig 1 ✓

Problematic assembly

contig 1
contig 2 ✗
contig 3
contig 4

# Haplotig purging

**HaploMerger2**

## HaploMerger2: rebuilding both haploid sub-assemblies from high-heterozygosity diploid genome assembly

Shengfeng Huang*, Mingjing Kang and Anlong Xu

## Identifying and removing haplotypic duplication in primary genome assemblies

Dengfeng Guan[1,2], Shane A. McCarthy [2], Jonathan Wood[3], Kerstin Howe [3], Yadong Wang[1,*] and Richard Durbin [2,3,*]
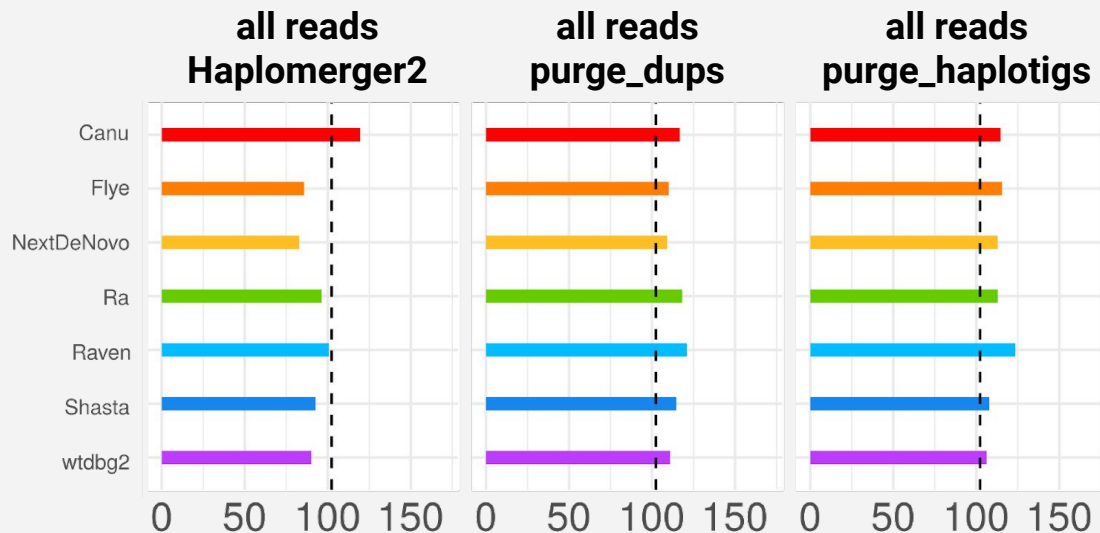
**purge_dups**

**Purge Haplotigs**

## Purge Haplotigs: allelic contig reassignment for third-gen diploid genome assemblies

Michael J. Roach* , Simon A. Schmidt and Anthony R. Borneman

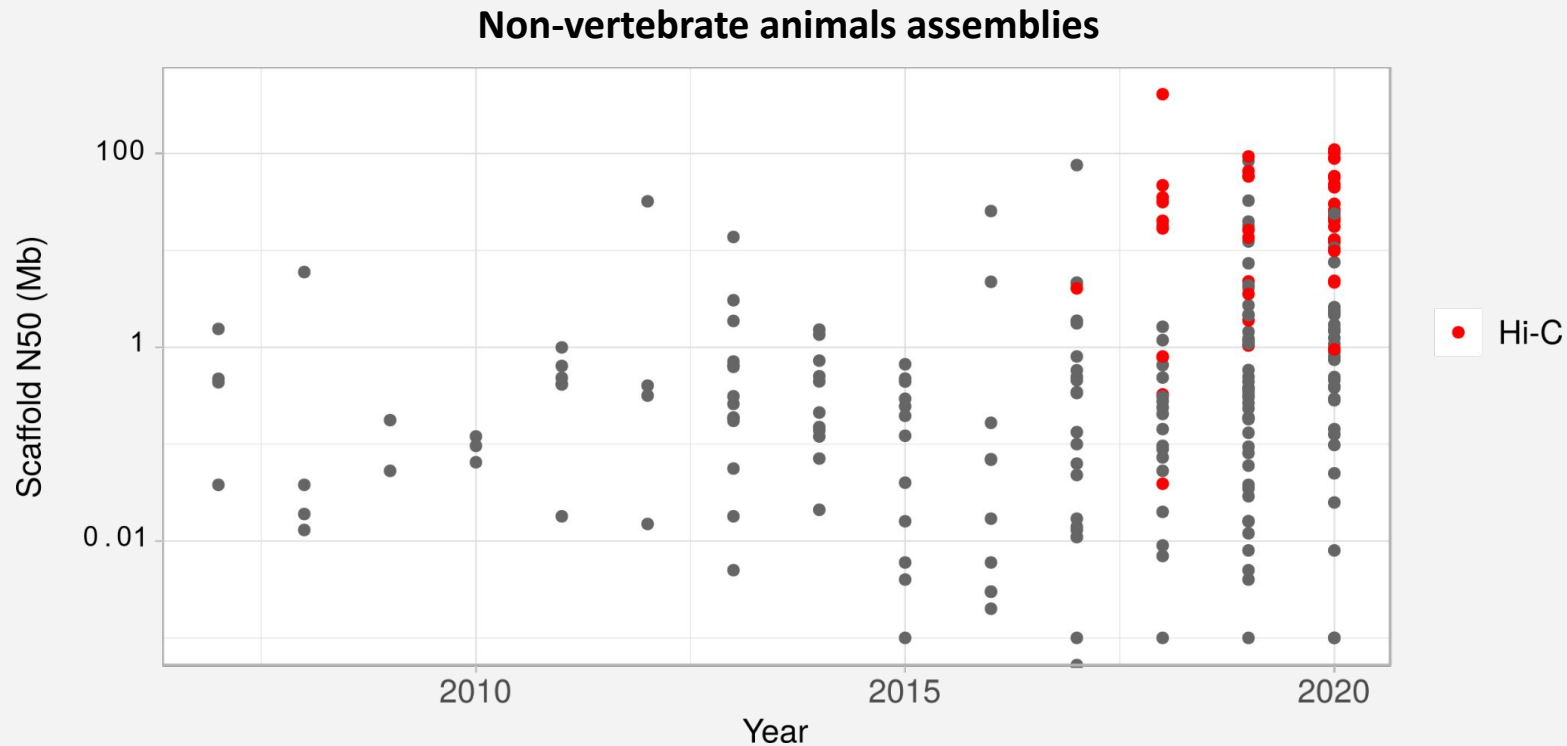# Haplotig purging

PacBio assemblies

# Scaffolding approaches

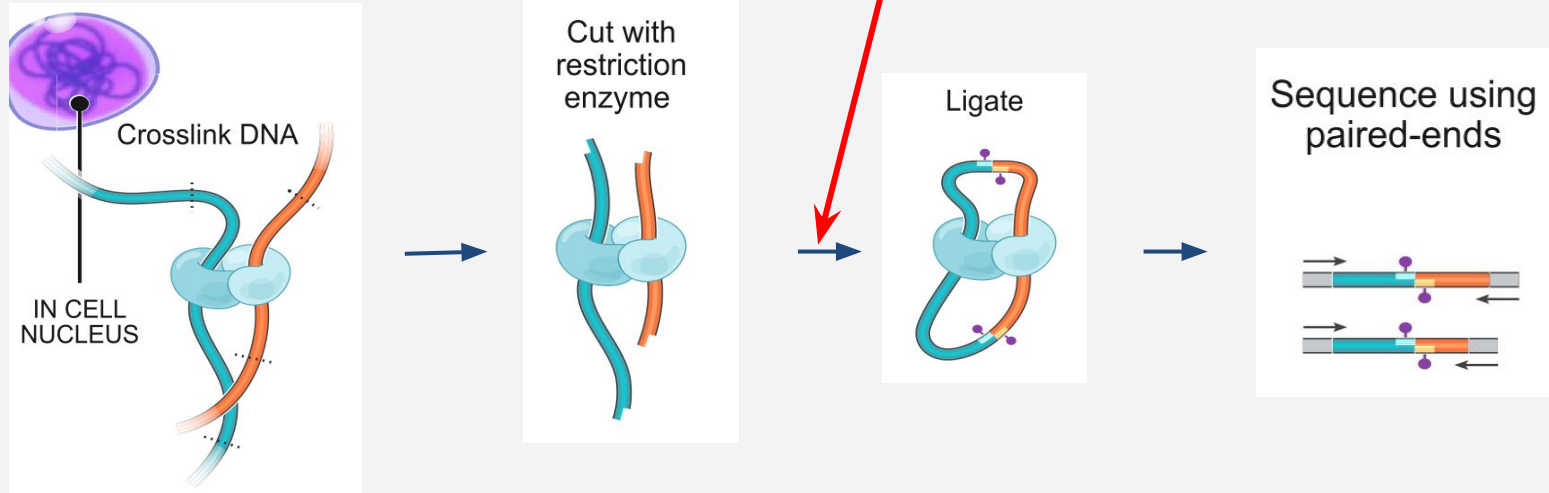Scaffolding: grouping and orienting contigs to build chromosome-level scaffolds

- ▸ **Long reads**

- ▸ **Linked reads:** barcoded short reads

- ▸ **Hi-C**

# **Scaffolding approaches:** Hi-C scaffolding



Non-vertebrate animals assemblies

# **Scaffolding approaches:** Hi-C scaffolding
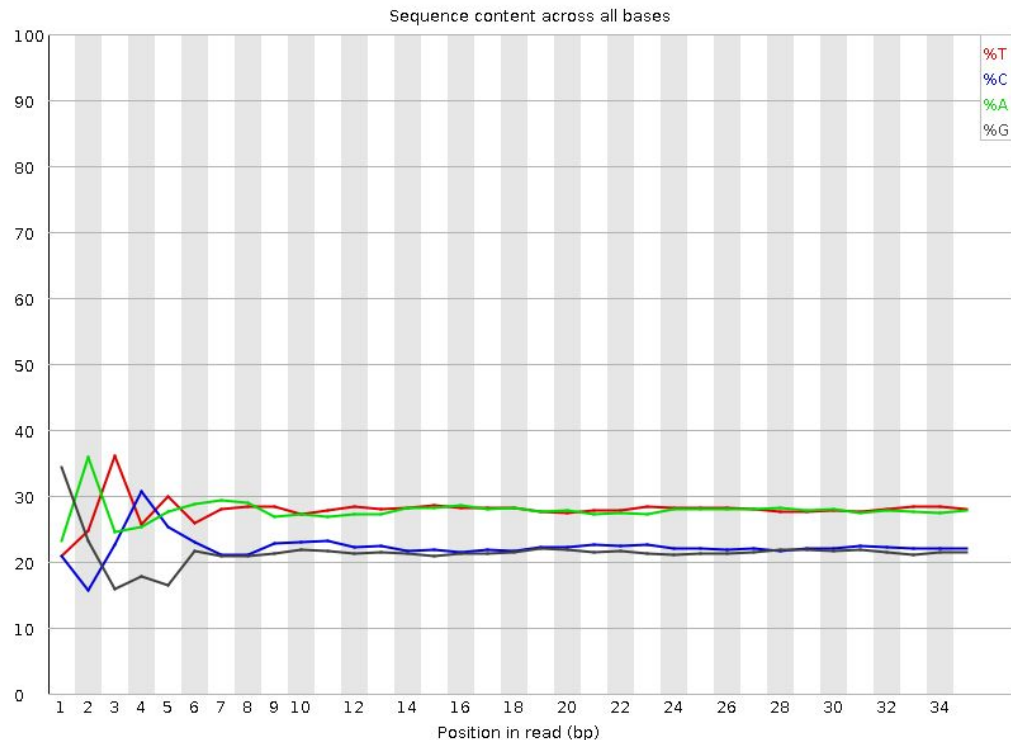
**Hi-C**

biotinylation



A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. Rao et *al.*, 2014
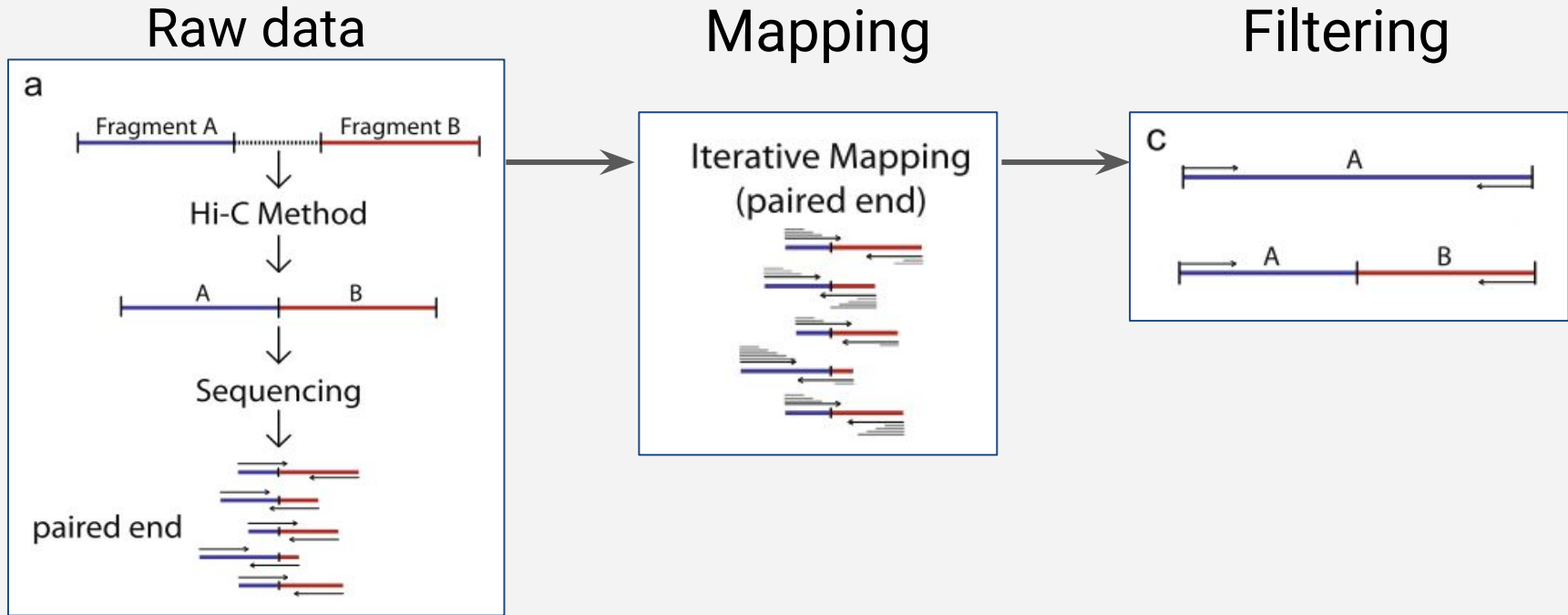
# **Scaffolding approaches:** Hi-C scaffolding

# **Scaffolding approaches:** Hi-C scaffolding

Raw data                    Mapping                    Filtering
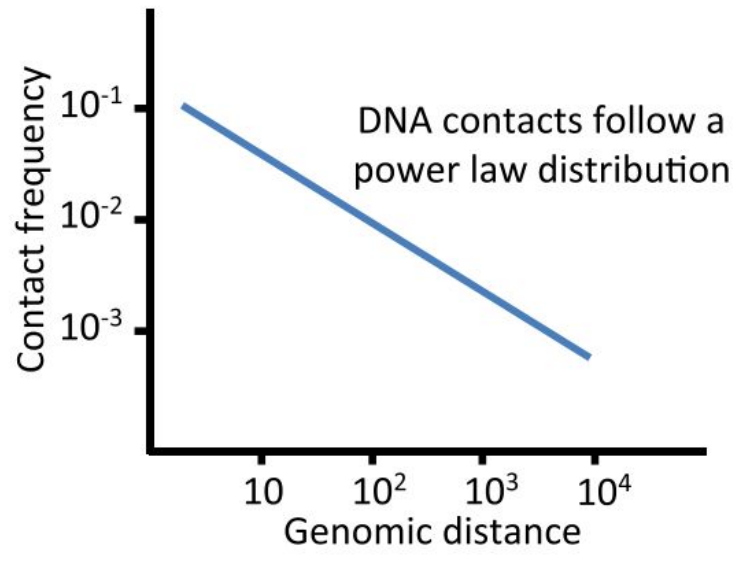


The Hitchhiker's guide to Hi-C analysis: Practical guidelines. Lajoie et *al.*, 2015

# **Scaffolding approaches:** Hi-C scaffolding



Contact map of
*Caenorhabditis elegans*

contact frequency = f(genomic distance)



DNA contacts follow a power law distribution

Contact genomics: scaffolding and phasing (meta)genomes using chromosome 3D physical signatures. Flot et *al.*, 2015

# **Scaffolding approaches:** Hi-C scaffolding

**High-throughput genome scaffolding from *in vivo* DNA interaction frequency**

Noam Kaplan & Job Dekker

**dnaTri**

**Lachesis**

**Chromosome-scale scaffolding of *de novo* genome assemblies based on chromatin interactions**

Joshua N Burton, Andrew Adey, Rupali P Patwardhan, Ruolan Qiu, Jacob O Kitzman & Jay Shendure

**High-quality genome (re)assembly using chromosomal contact data**

Hervé Marie-Nelly, Martial Marbouty, Axel Cournac, Jean-François Flot, Gianni Liti, Dante Poggi Parodi, Sylvie Syan, Nancy Guillén, Antoine Margeot, Christophe Zimmer & Romain Koszul

**GRAAL**

# Scaffolding approaches: Hi-C scaffolding

De novo assembly of the *Aedes aegypti* genome using Hi-C yields chromosome-length scaffolds

Olga Dudchenko[1,2,3,4], Sanjit S. Batra[1,2,3,*], Arina D. Omer[1,2,3,*], Sarah K. Nyquist[1,3], Marie Hoeger[1,3], Neva C. Durand[1,...]

**3D-DNA**

**SALSA2**

Integrating Hi-C links with assembly graphs for chromosome-scale assembly

Jay Ghurye, Arang Rhie, Brian P. Walenz, Anthony Schmitt, Siddarth Selvaraj, Mihai Pop, Adam M. Phillippy ✉, Sergey Koren ✉

instaGRAAL: chromosome-level quality scaffolding of genomes using a proximity ligation-based scaffolder

Lyam Baudry, Nadège Guiglielmoni, Hervé Marie-Nelly, Alexandre Cormier, Martial Marbouty, Komlan Avia, Yann Loe Mie, Olivier Godfroy, Lieven Sterck, J. Mark Cock, Christophe Zimmer, Susana M. Coelho ✉ & Romain Koszul ✉

**instaGRAAL**

# **Scaffolding approaches:** Hi-C scaffolding

And in 2021

**EndHiC: assemble large contigs into chromosomal-level scaffolds using the Hi-C links from contig ends**

Sen Wang, Hengchao Wang, Fan Jiang, Anqi Wang, Hangwei Liu, Hanbo Zhao, Boyuan Yang, Dong Xu, Yan Zhang, Wei Fan

## Efficient iterative Hi-C scaffolder based on N-best neighbors

Dengfeng Guan[1,2,4], Shane A. McCarthy[2,3], Zemin Ning[3], Guohua Wang[1*], Yadong Wang[1*] and Richard Durbin[2,3*]

## YaHS: yet another Hi-C scaffolding tool

Chenxi Zhou[1], Shane A. McCarthy[1,2], Richard Durbin[1,2]

# **Scaffolding approaches:** Hi-C scaffolding

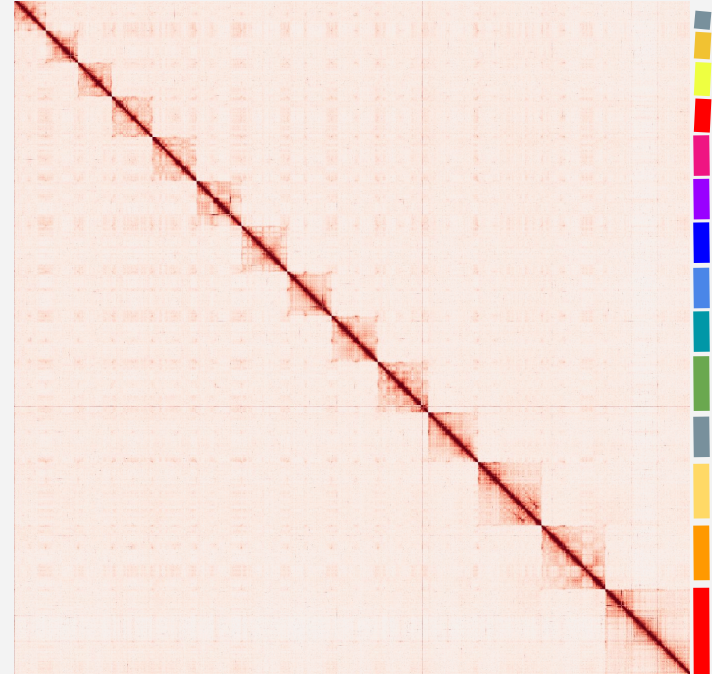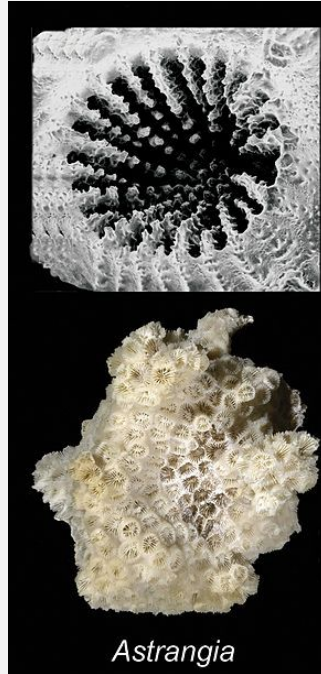

www.dnazoo.org

# **Scaffolding approaches:** Hi-C scaffolding

*Astrangia poculata*

(coral)

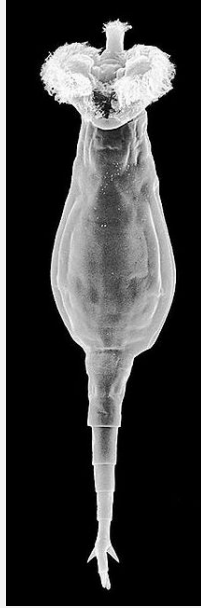14 scaffolds

455 Mb





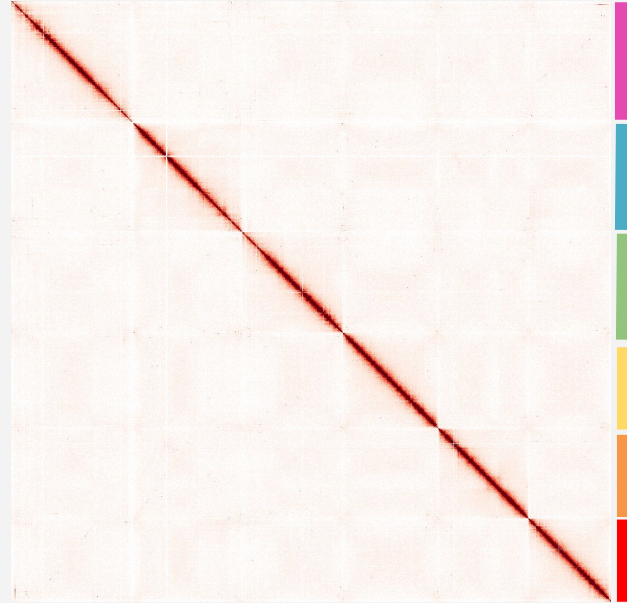**Hi-C contact map of *Astrangia poculata***

# Scaffolding approaches: Hi-C scaffolding

*Adineta vaga* (rotifer)

6 scaffolds



Who Needs Sex (or Males) Anyway?
Liza Gross, PloS Biology, 2007

**Hi-C contact map of *Adineta vaga***

# **Scaffolding approaches:** Hi-C scaffolding

"What coverage should I get?"

➔ Arima recommends 200 millions pairs per Gb

| Species | Size | # fragments | # Hi-C pairs | Hi-C mapping |
|---|---|---|---|---|
| *Adineta vaga* | 101 Mb | 30 | 55 millions | 83% |
| *Astrangia poculata* | 455 Mb | 2995 | 723 millions | 67% |
| *Flaccisagitta enflata* | 929 Mb | 6612 | 489 millions | 37% |
| *Mercenaria mercenaria* | 1.86 Gb | 5118 | 455 millions | 55% |

# And then...

▸ **Gap filling**: TGS-GapCloser...

▸ **Polishing**: using high-accuracy reads, HyPo, Racon...

# Gap filling & Polishing

|  | Scaffolds | After TGS-Gapcloser | After HyPo |
|---|---|---|---|
| *Flaccisagitta enflata* | 9,239 | 3,694 | 1,476 |
| *Norana najaformis* | 860 | 748 | 632 |
| *Lucinoma borealis* | 24,786 | 5,093 | 2,135 |

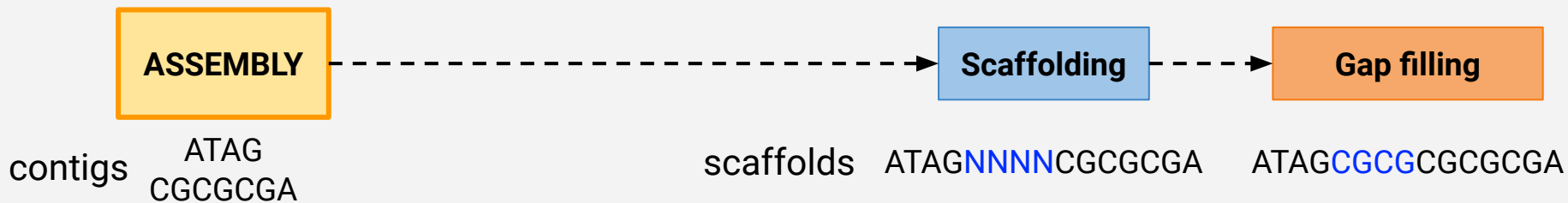# Assembly pipeline

**ASSEMBLY**

reads    ATTTGTACG
            GTACGGACA
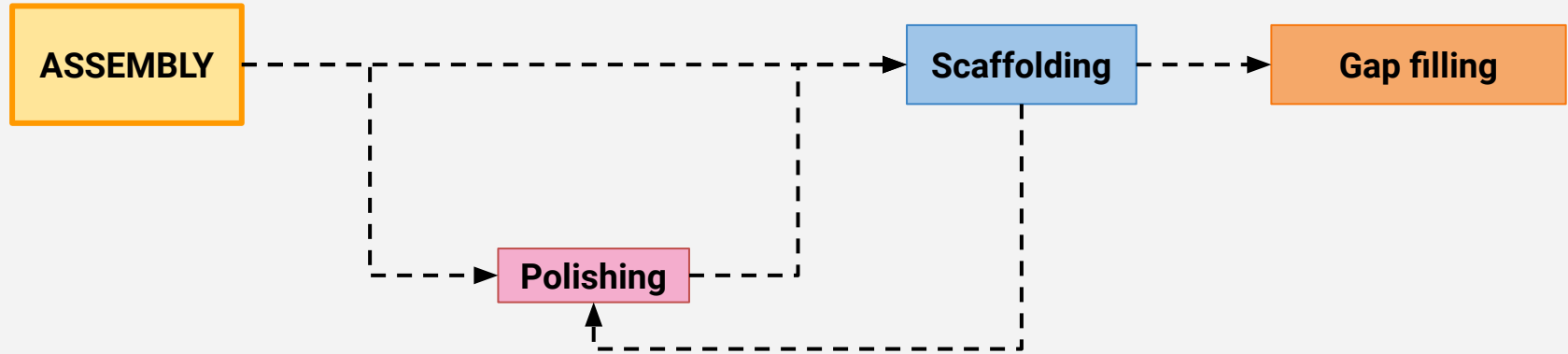                 GGACATAGTA

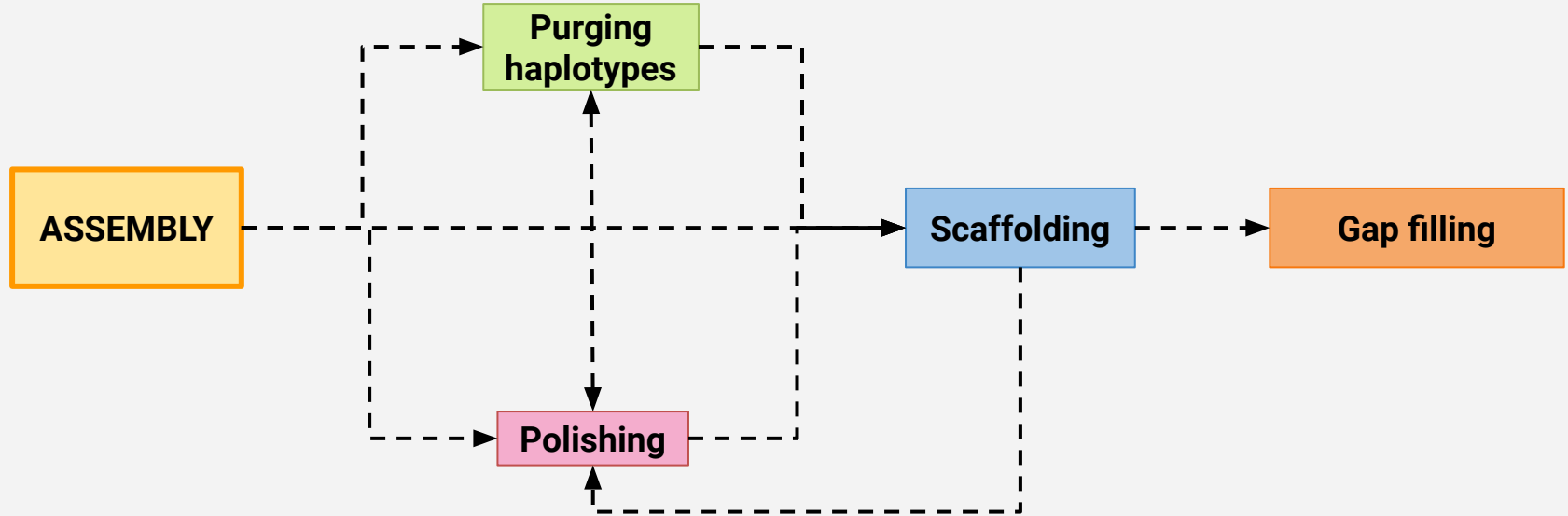contig    ATTTGTACGGACATAGTA

# Assembly pipeline

**ASSEMBLY**

contigs  ATAG
CGCGCGA

**Scaffolding**

scaffolds  ATAGNNNNCGCGCGA

# Assembly pipeline



contigs    ATAG CGCGCGA

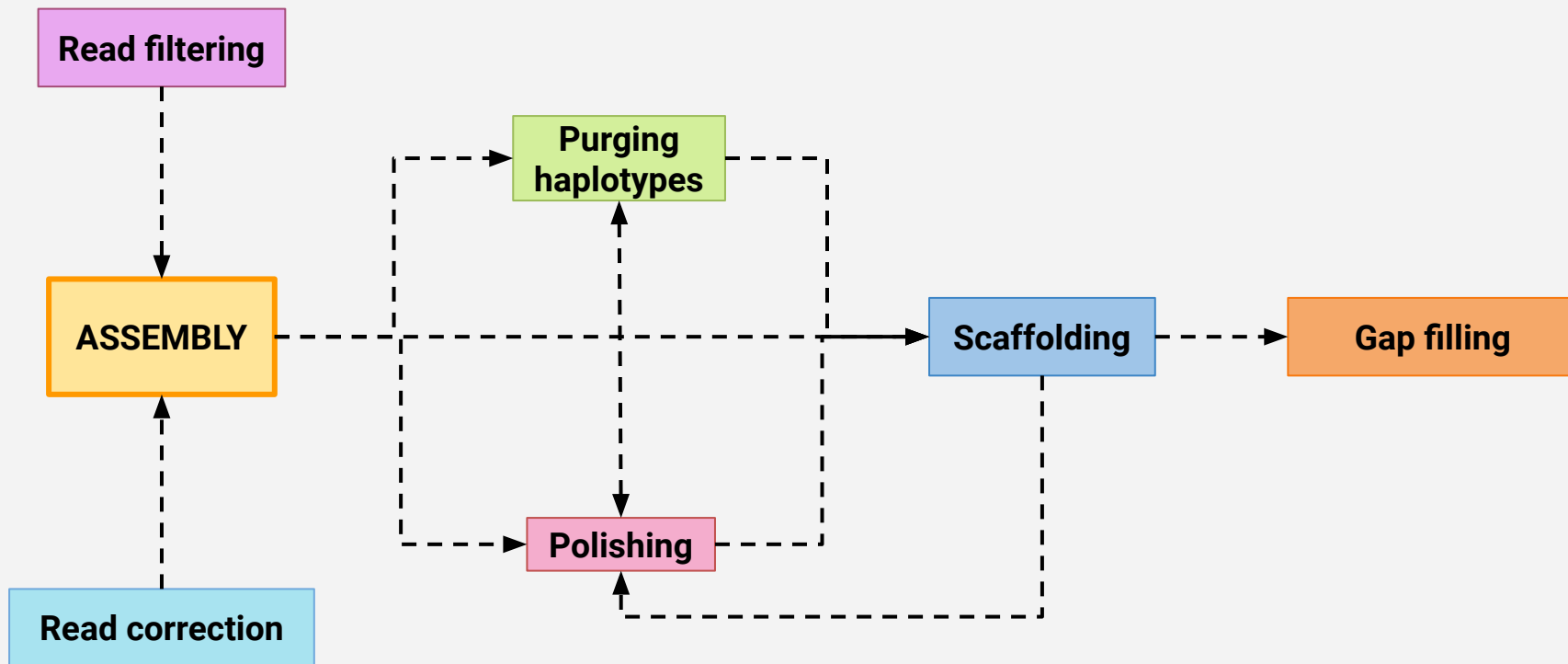scaffolds    ATAGNNNNCGCGCGA

ATAGCGCGCGCGCGA

# Assembly pipeline

# Assembly pipeline

# Assembly pipeline

# Assembly pipeline