# Genome assembly

**Meleshko Dmitry**
*meleshko.dmitrii@gmail.com*

# *De novo* genome assembly
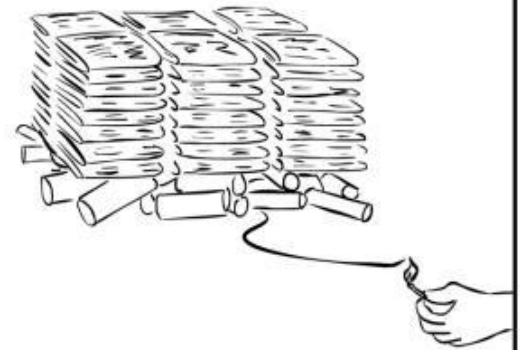
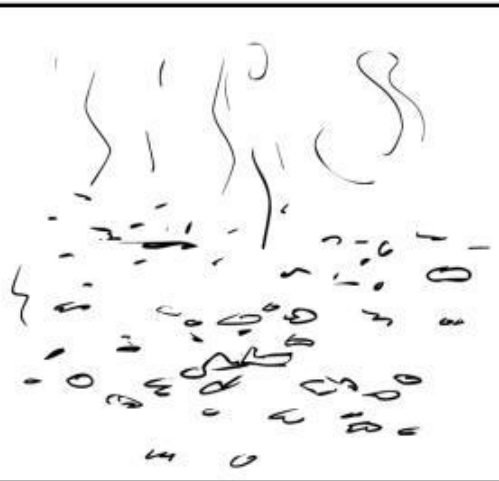# *De novo* whole genome assembly

# Whole genome assembly

# Shortest common supersequence

Given a set of strings $\{s_1, \ldots, s_n\}$, find a shortest string $S$ containing each $s_i$ as a substring

# Shortest common supersequence

Given a set of strings $\{s_1, \ldots, s_n\}$, find a shortest string $S$ containing each $s_i$ as a substring

Is NP-complete

Has nothing to do with real genome assembly problem

# Why to assemble?

- NGS
  - Billions of short reads
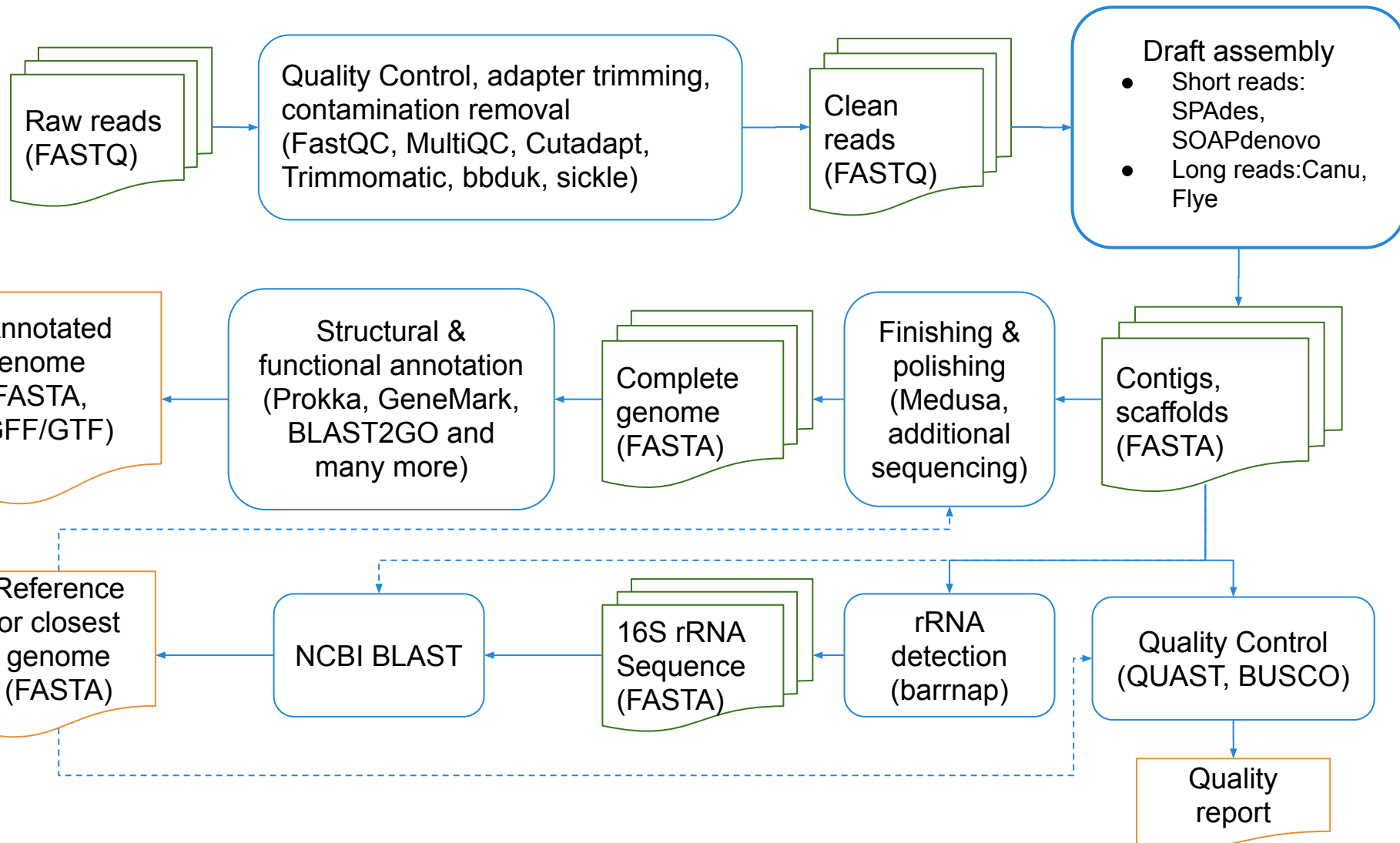  - Sequencing errors
  - Contaminants

- Assembly
  - ✓ Corrects sequencing errors
  - ✓ Much longer sequences
  - ✓ Each genomic region is presented only once
  - ✗ May introduce errors
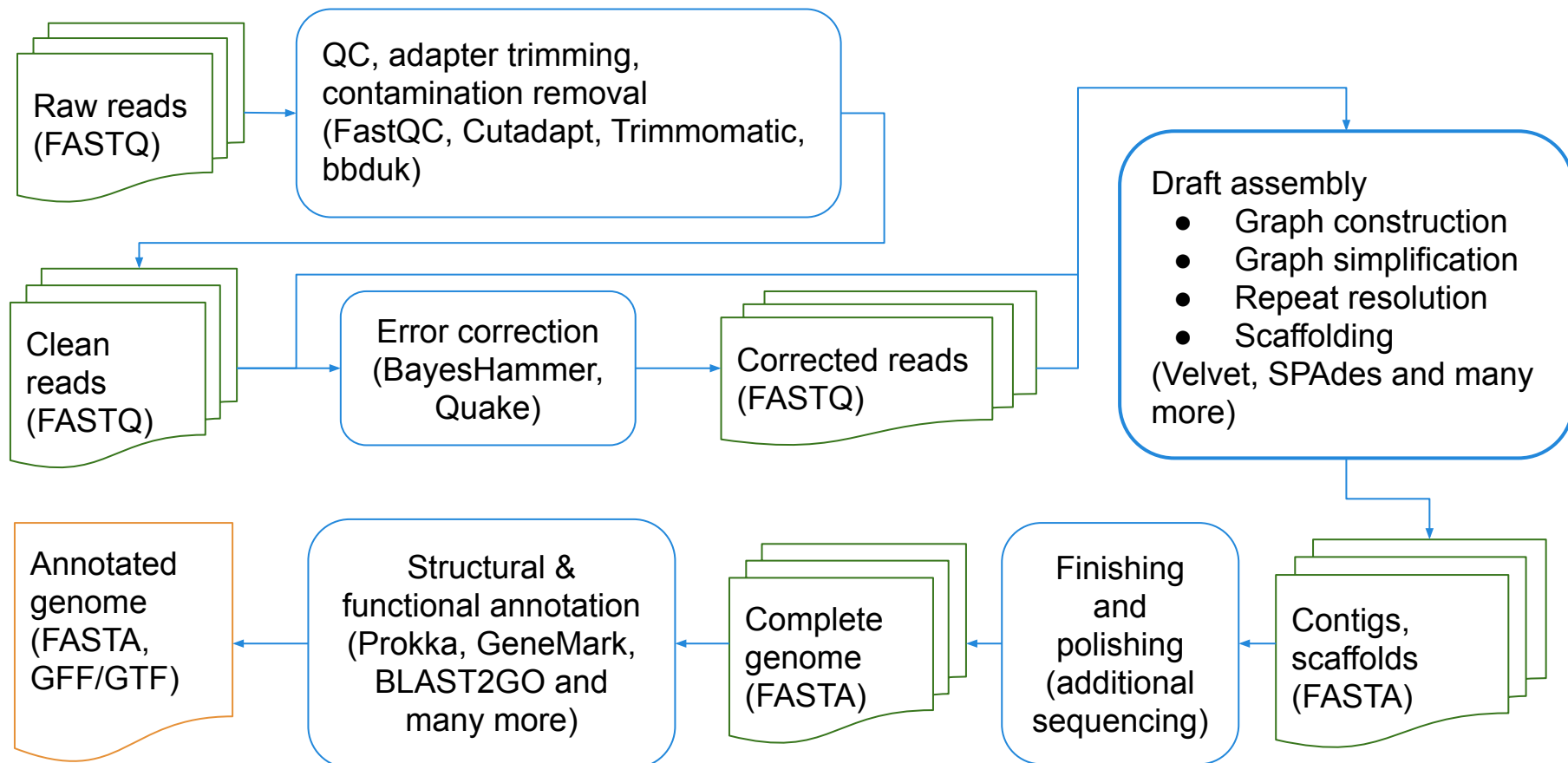
Hard to perform analysis

# Assembly types

- *De novo* genome assembly
  - Long reads
  - Short reads
  - Hybrid
- Reference-assisted genome assembly
  - Closely related species
- Transcriptome assembly
  - *De novo*
  - Reference based

# *De novo* genome assembly

Raw reads (FASTQ) → Quality Control, adapter trimming, contamination removal (FastQC, MultiQC, Cutadapt, Trimmomatic, bbduk, sickle) → Clean reads (FASTQ) → Draft assembly
- Short reads: SPAdes, SOAPdenovo
- Long reads: Canu, Flye

Annotated genome (FASTA, GFF/GTF) ← Structural & functional annotation (Prokka, GeneMark, BLAST2GO and many more) ← Complete genome (FASTA) ← Finishing & polishing (Medusa, additional sequencing) ← Contigs, scaffolds (FASTA)

Reference or closest genome (FASTA) ← NCBI BLAST ← 16S rRNA Sequence (FASTA) ← rRNA detection (barrnap) ← Quality Control (QUAST, BUSCO) → Quality report

# *De novo* genome assembly

Raw reads (FASTQ) → QC, adapter trimming, contamination removal (FastQC, Cutadapt, Trimmomatic, bbduk) → Clean reads (FASTQ) → Error correction (BayesHammer, Quake) → Corrected reads (FASTQ) → Draft assembly
- Graph construction
- Graph simplification
- Repeat resolution
- Scaffolding

(Velvet, SPAdes and many more) → Contigs, scaffolds (FASTA) → Finishing and polishing (additional sequencing) → Complete genome (FASTA) → Structural & functional annotation (Prokka, GeneMark, BLAST2GO and many more) → Annotated genome (FASTA, GFF/GTF)

# Assembling Sanger reads

# Early days

- Sanger sequencing
  - Long reads
  - Low coverage

- Overlap-Layout-Consensus (OLC)
  - Find overlaps between all reads
  - Order reads
  - Merge into consensus sequence

# Finding overlaps

- Align reads all-to-all

  - BLAST and similar algorithms

- Ignore "insufficient" overlaps

# Finding overlaps

- Align reads all-to-all

  ○ BLAST and similar algorithms

- Ignore "insufficient" overlaps

  ○ At least 40bp

  ○ >94% similarity

# Assembly example

# Assembly example

# Assembly example

# Overlap graph

# Overlap graph

# Layout

# Layout

# Layout

# Layout



A → B → C → E → D

# Layout

# Consensus

# NGS and OLC

- Overlap-Layout-Consensus is not applicable

  - Hard to find overlaps between short reads

  - Impossible to scale to such amount of reads

- De Bruijn graph approach (Pevzner et al., 2001)

- String Graph approach (Meyers, 2005)

# NGS era

# De Bruijn graph in a nutshell

**He that mischief hatches, mischief catches**

Sequencing

**He that mischief**

**mischief hatches,**

**hatches, mischief**

**, mischief catches**

# De Bruijn graph in a nutshell

**, mischief catches**

**mischief hatches,**

**He that mischief**

**hatches, mischief**

# De Bruijn graph in a nutshell

# De Bruijn graph

**ACGTCCGTAA**

**<span style="color:red">AC</span>GTCCGTAA**

k=2

# De Bruijn graph

**ACGTCCGTAA**

k=2

( AC )     ( CG )

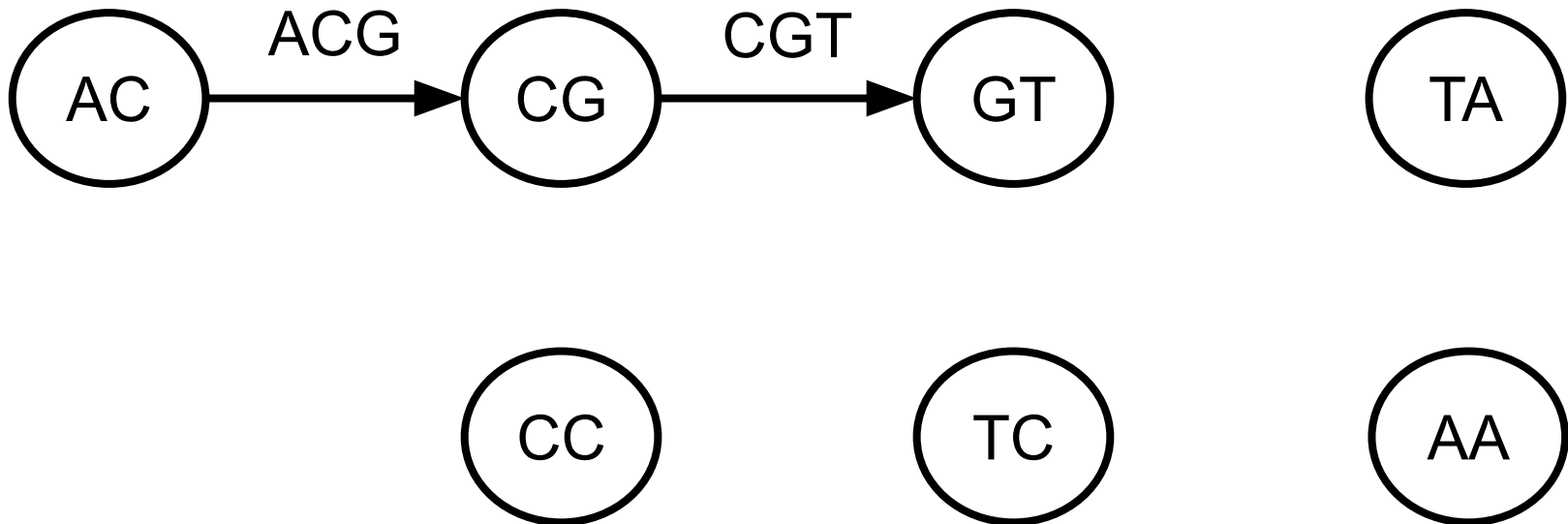# De Bruijn graph

**AC<span style="color:red">GT</span>CCGTAA**

k=2

# De Bruijn graph

**ACG<span style="color:red">TC</span>CGTAA**

k=2

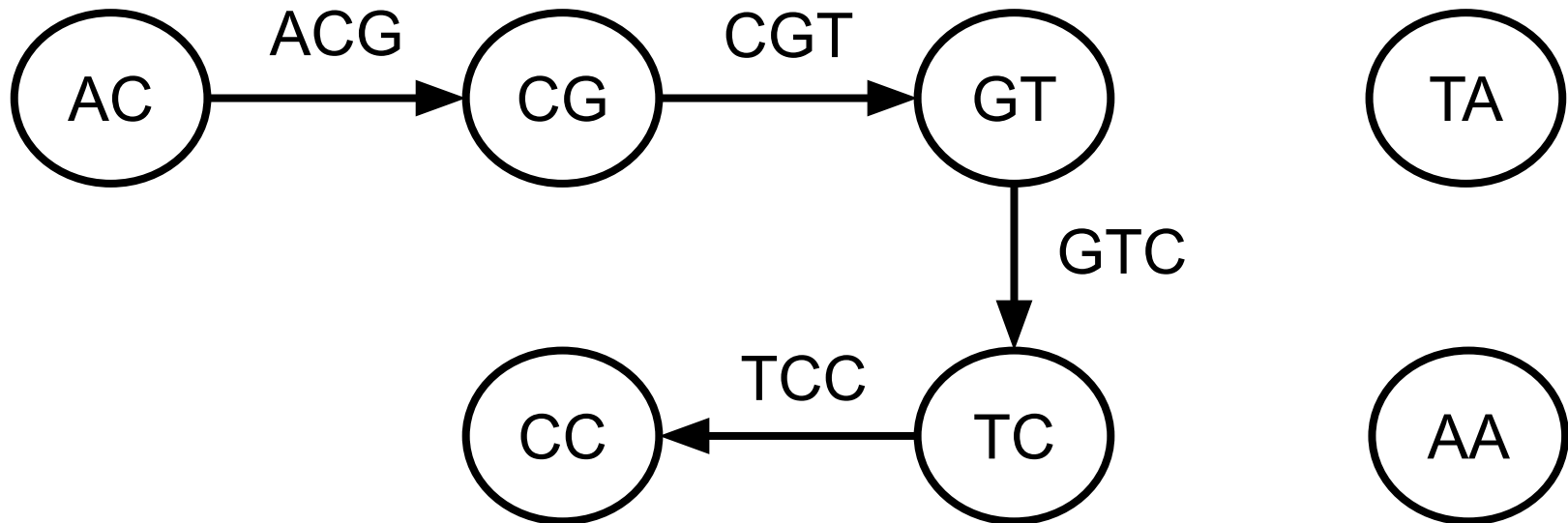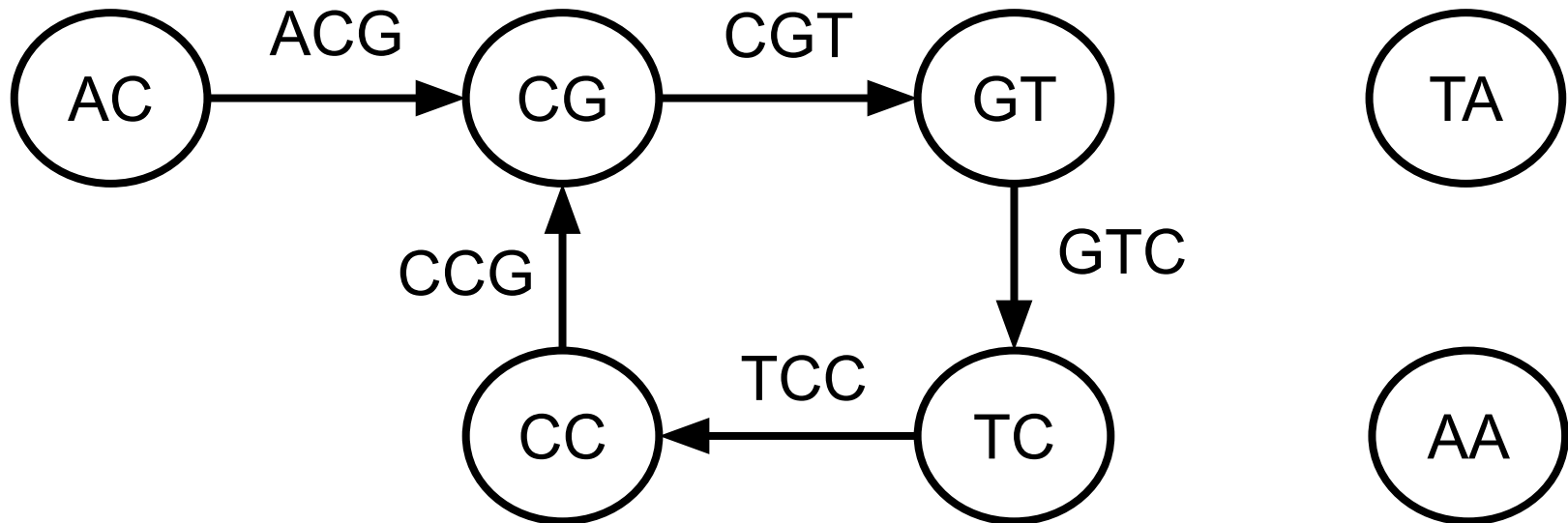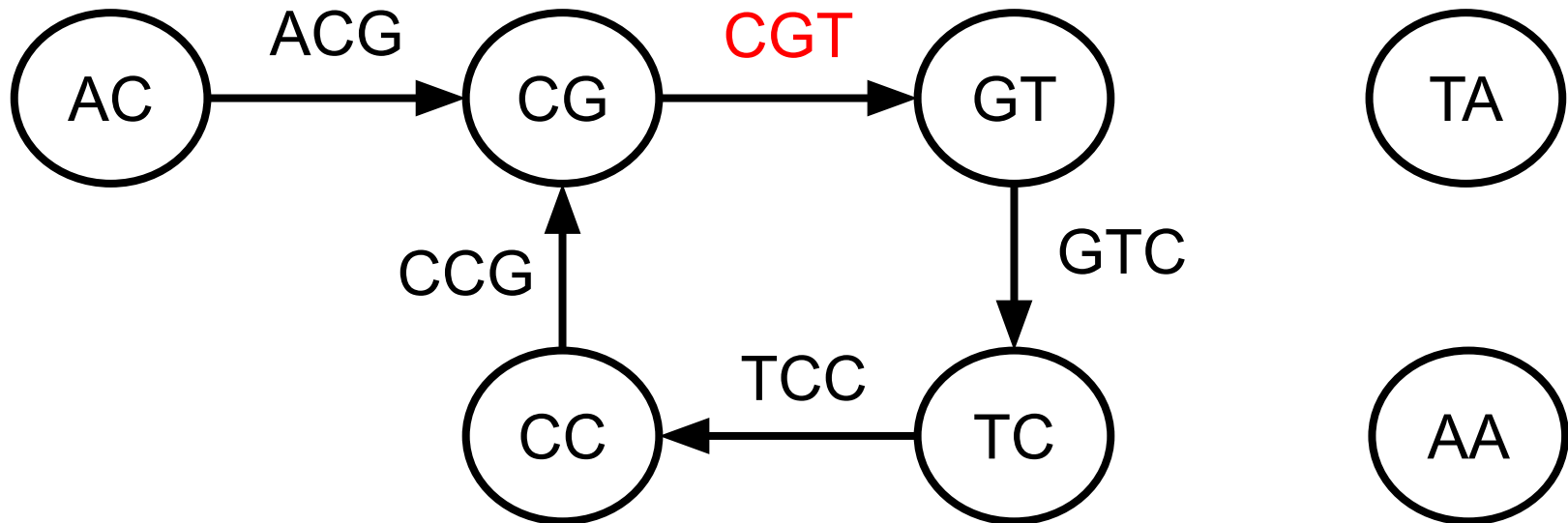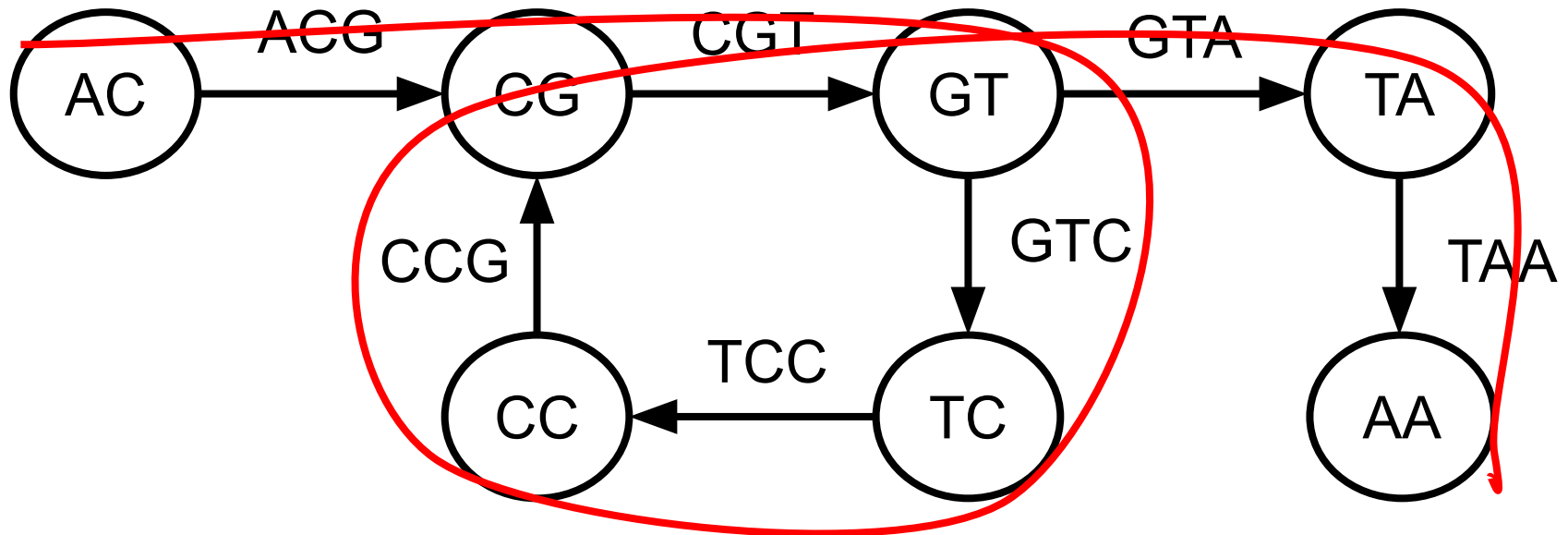# De Bruijn graph

**ACGT<span style="color:red">CC</span>GTAA**

k=2

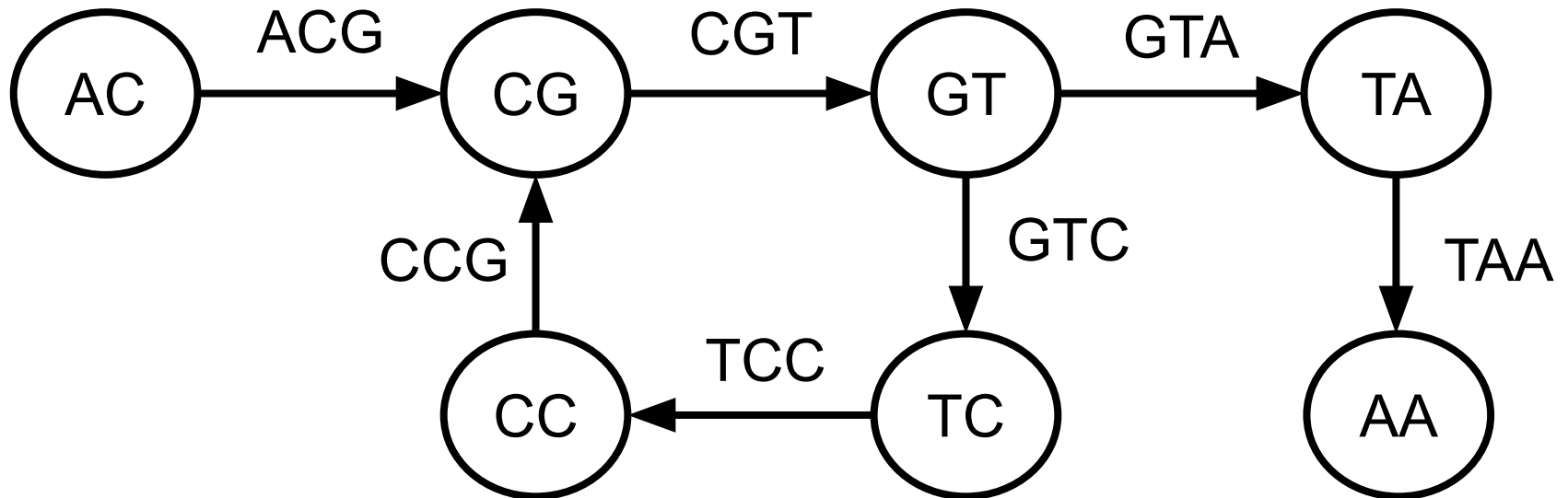# De Bruijn graph

**ACGTC<span style="color:red">CG</span>TAA**

k=2

# De Bruijn graph

**ACGTCC<span style="color:red">GT</span>AA**

k=2

AC

CG

GT

CC

TC

**ACGTCCG<span style="color:red">TA</span>A**

k=2

# De Bruijn graph

**ACGTCCGT**<span style="color:red">**AA**</span>

k=2

AC    CG    GT    TA

CC    TC    AA

# De Bruijn graph

**ACG**TCCGTAA

k=2

**A<span style="color:red">CGT</span>CCGTAA**

k=2

**AC<span style="color:red">GTC</span>CGTAA**

k=2

# ACG**TCC**GTAA

k=2

**ACGT<span style="color:red">CCG</span>TAA**

k=2

**ACGTC<span style="color:red">CGT</span>AA**

k=2
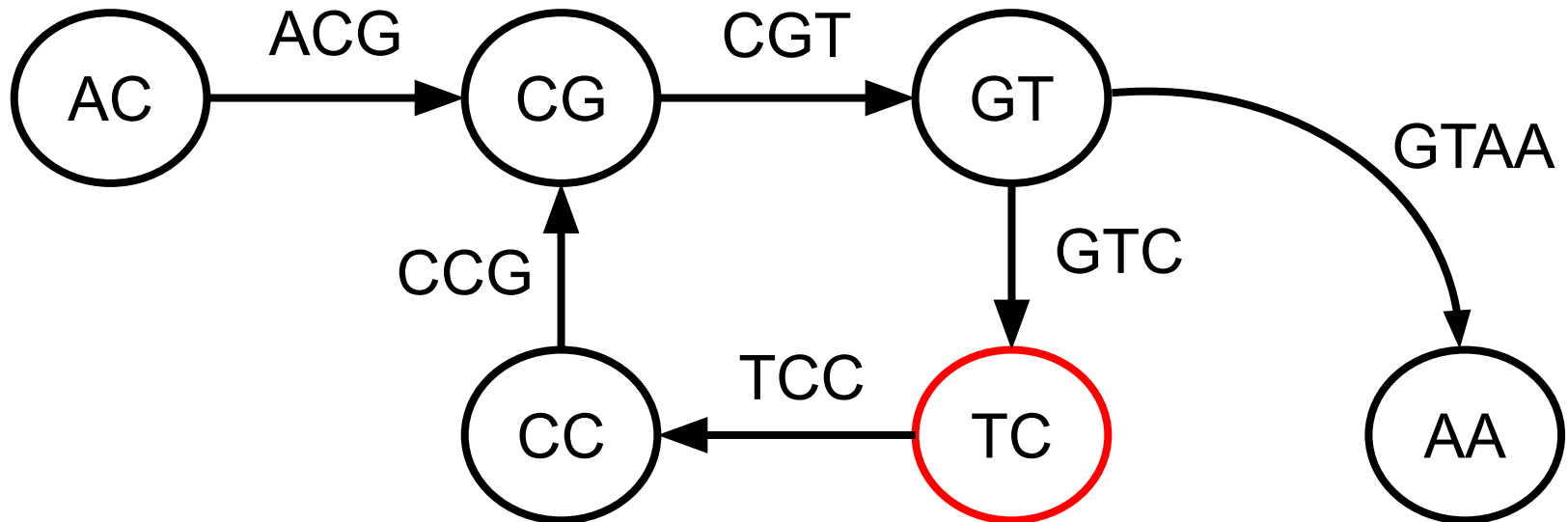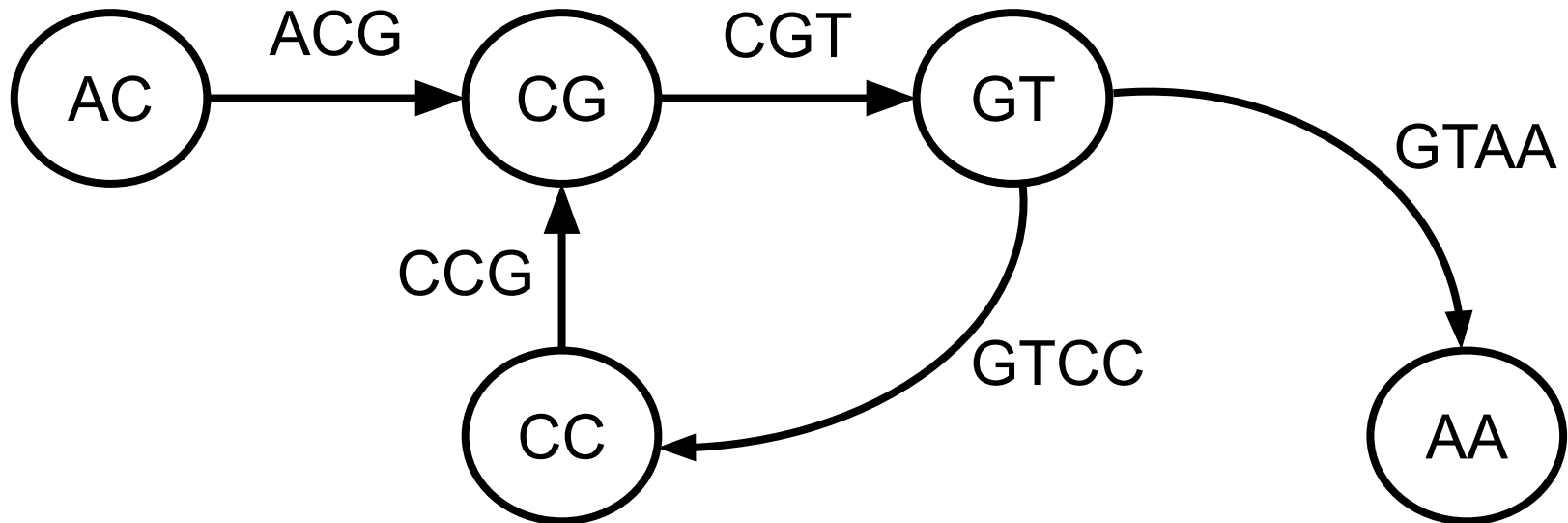
# ACGTCCGTAA

k=2

# ACGTCCGTAA

k=2

# Condensed de Bruijn graph

**ACGTCCGTAA**

k=2

# Condensed de Bruijn graph

**ACGTCCGTAA**

k=2

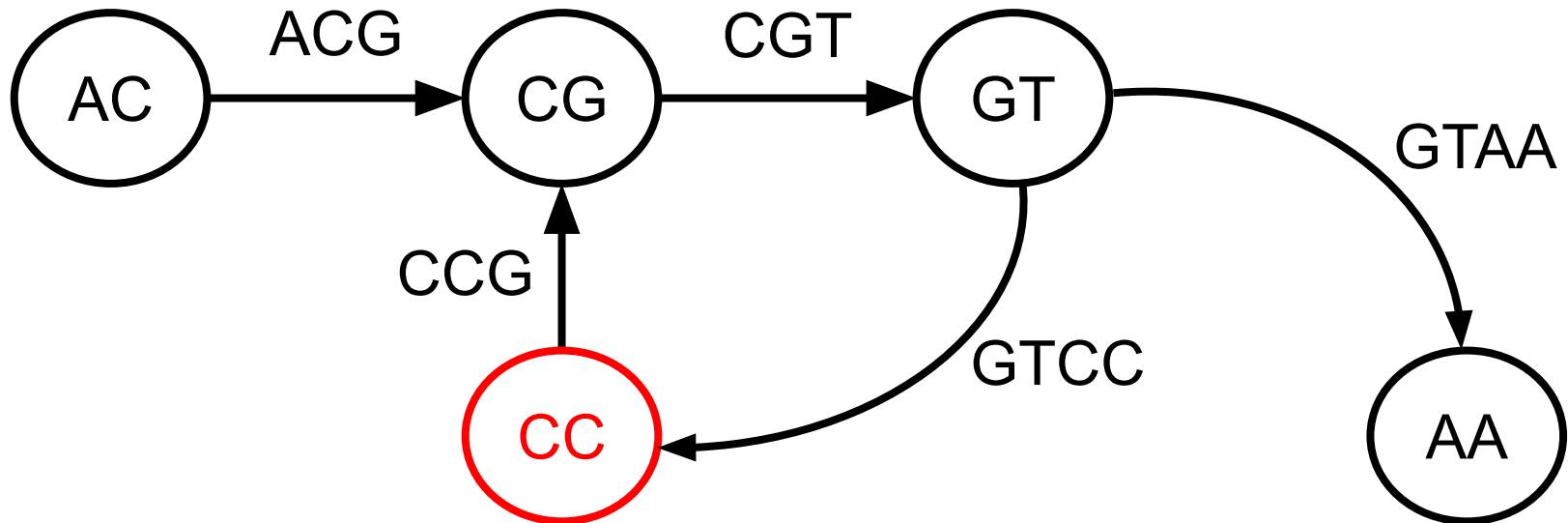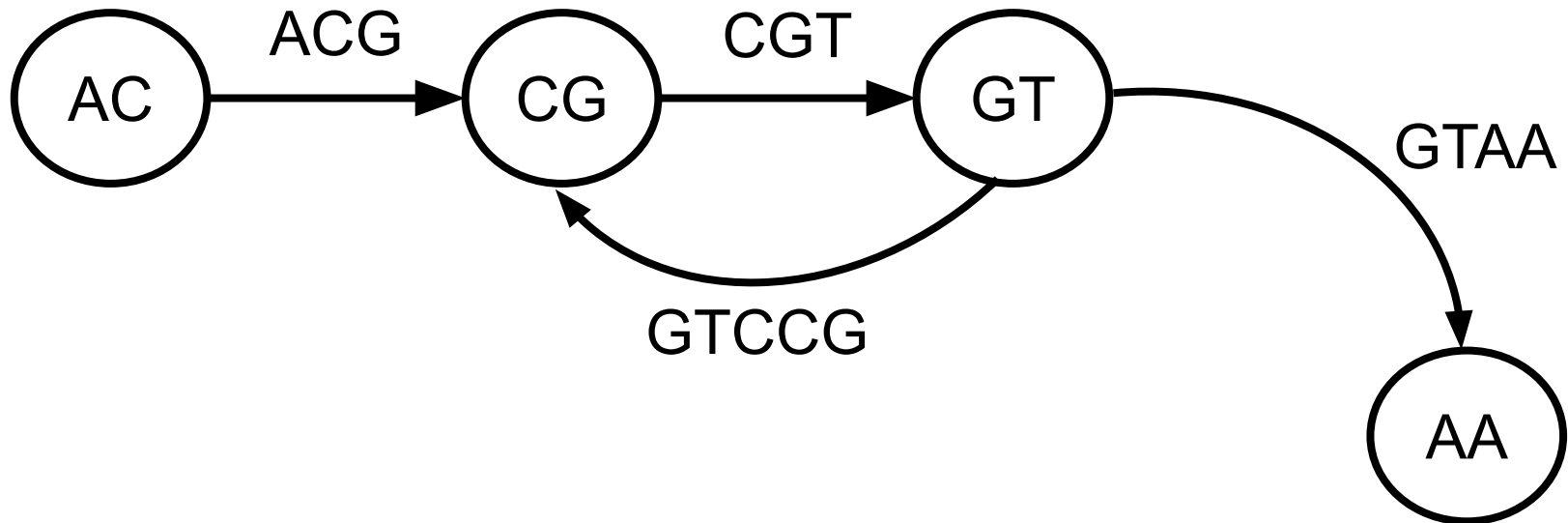# Condensed de Bruijn graph

**ACGTCCGTAA**

k=2

# Condensed de Bruijn graph

**ACGTCCGTAA**

k=2

# Condensed de Bruijn graph

ACGTCCGTAA

k=2

# Condensed de Bruijn graph

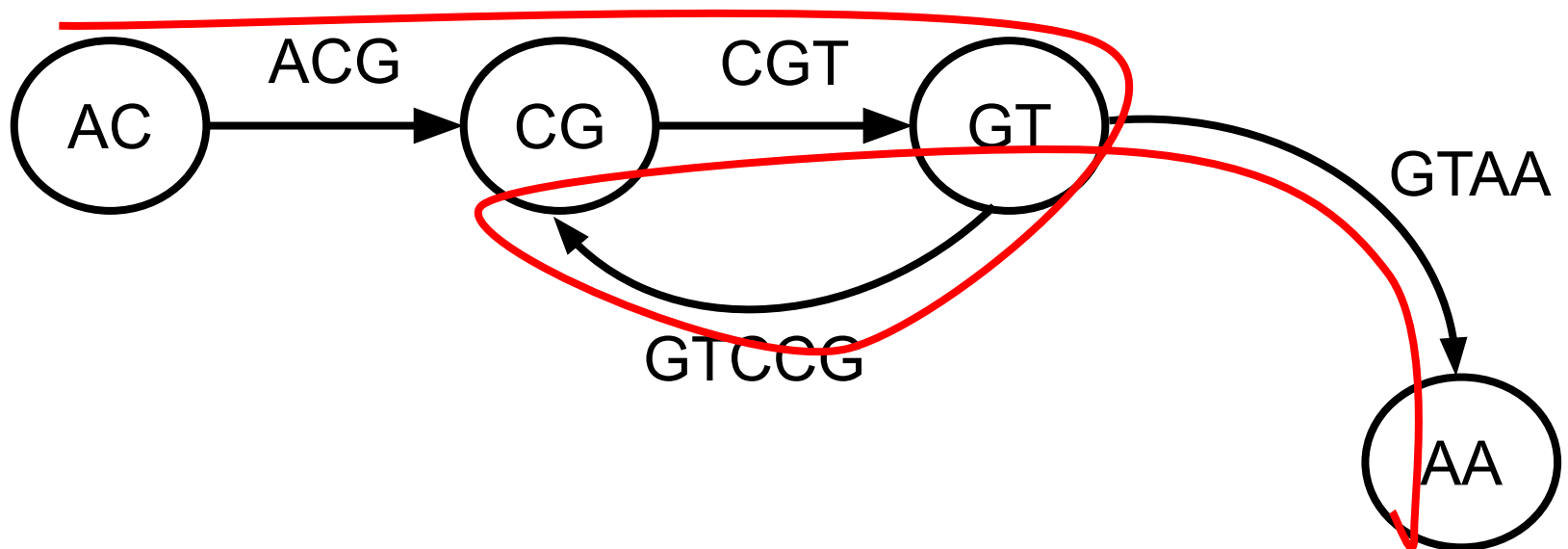**ACGTCCGTAA**

k=2

# Condensed de Bruijn graph

**ACGTCCGTAA**
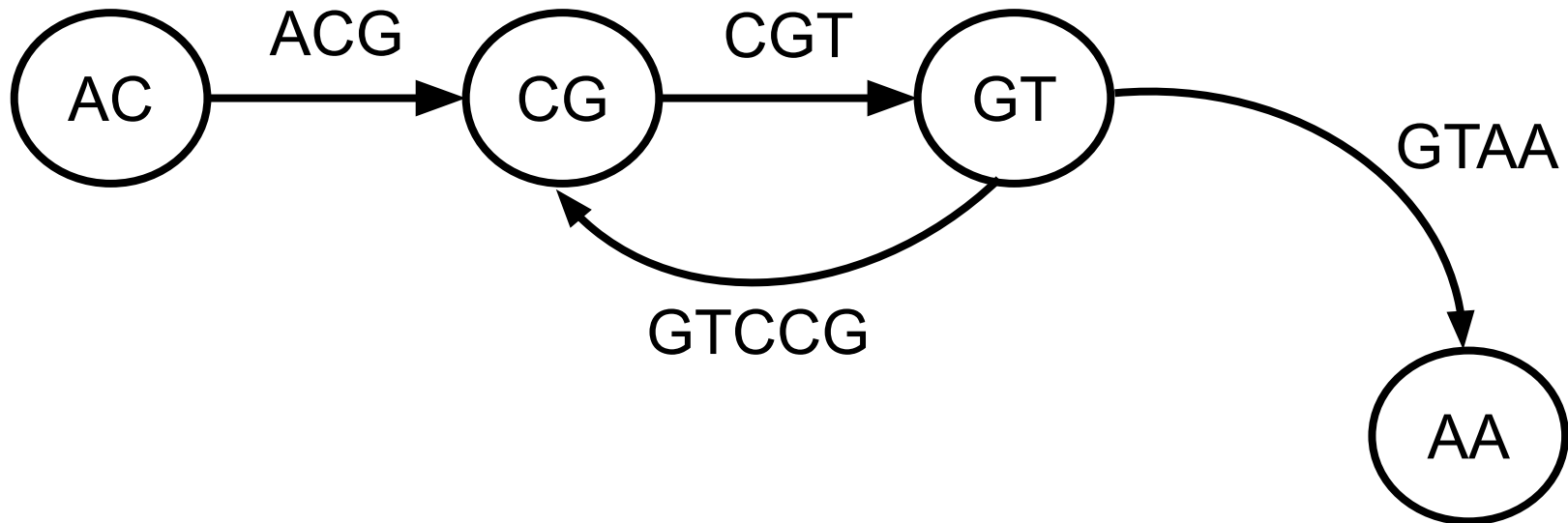
k=2

**ACGTCCGTAA**

k=2

# Repeats in de Bruijn graph

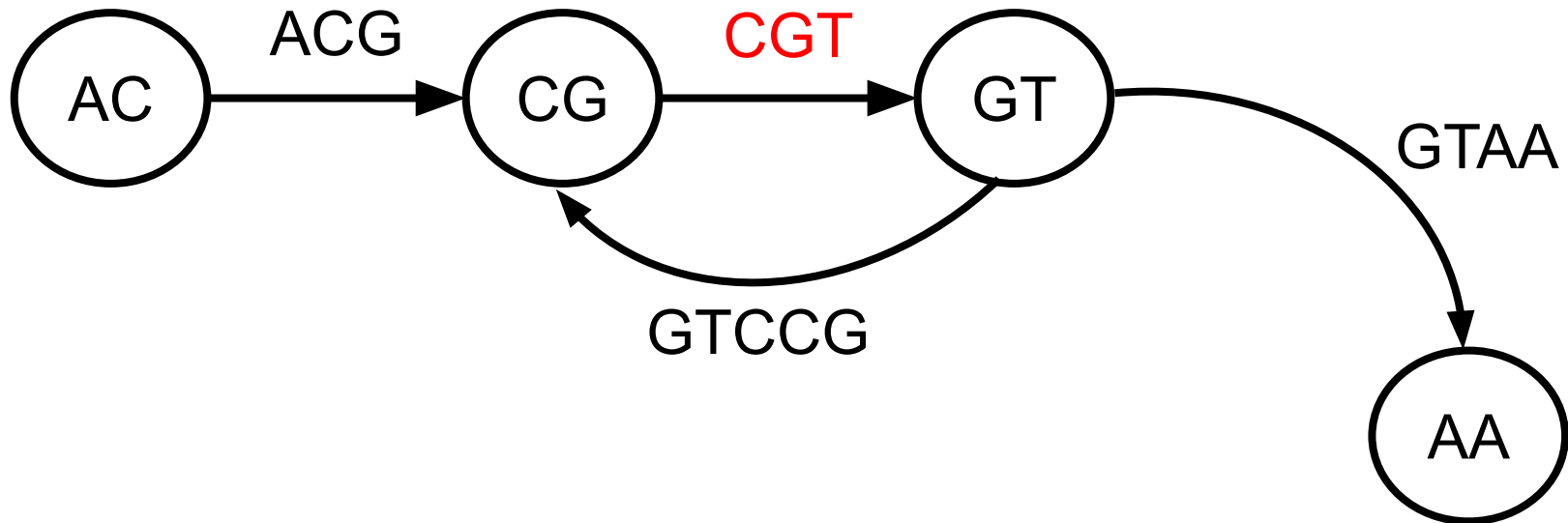**A<span style="color:red">CGT</span>C<span style="color:red">CGT</span>AA**
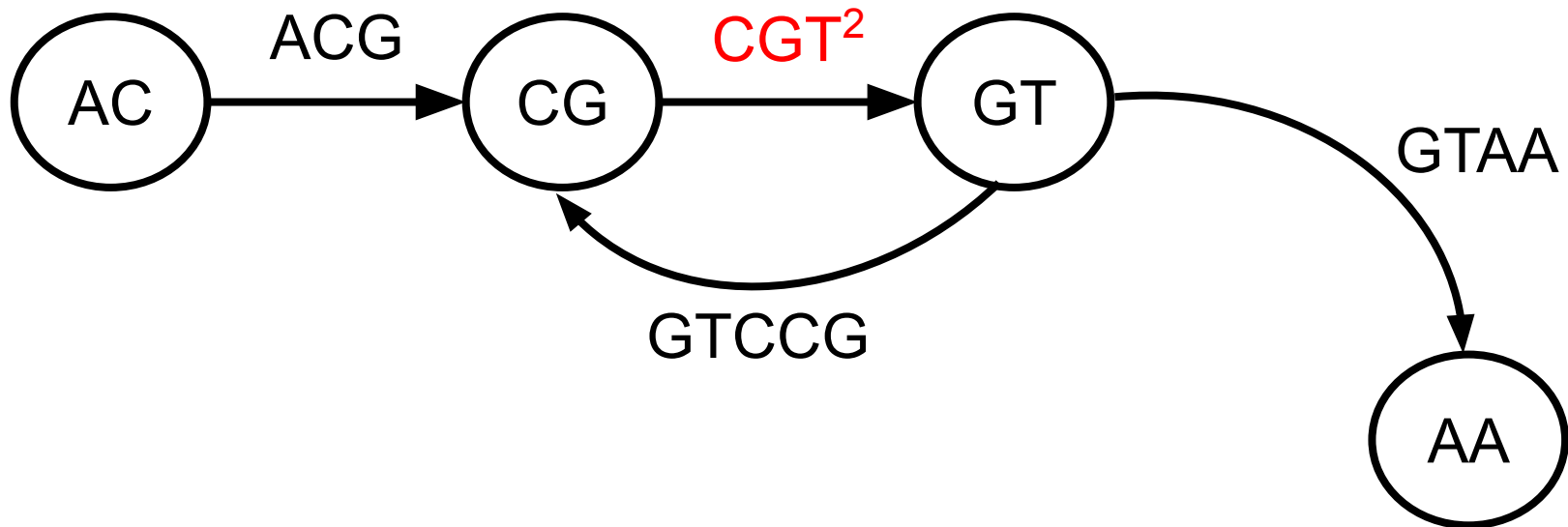
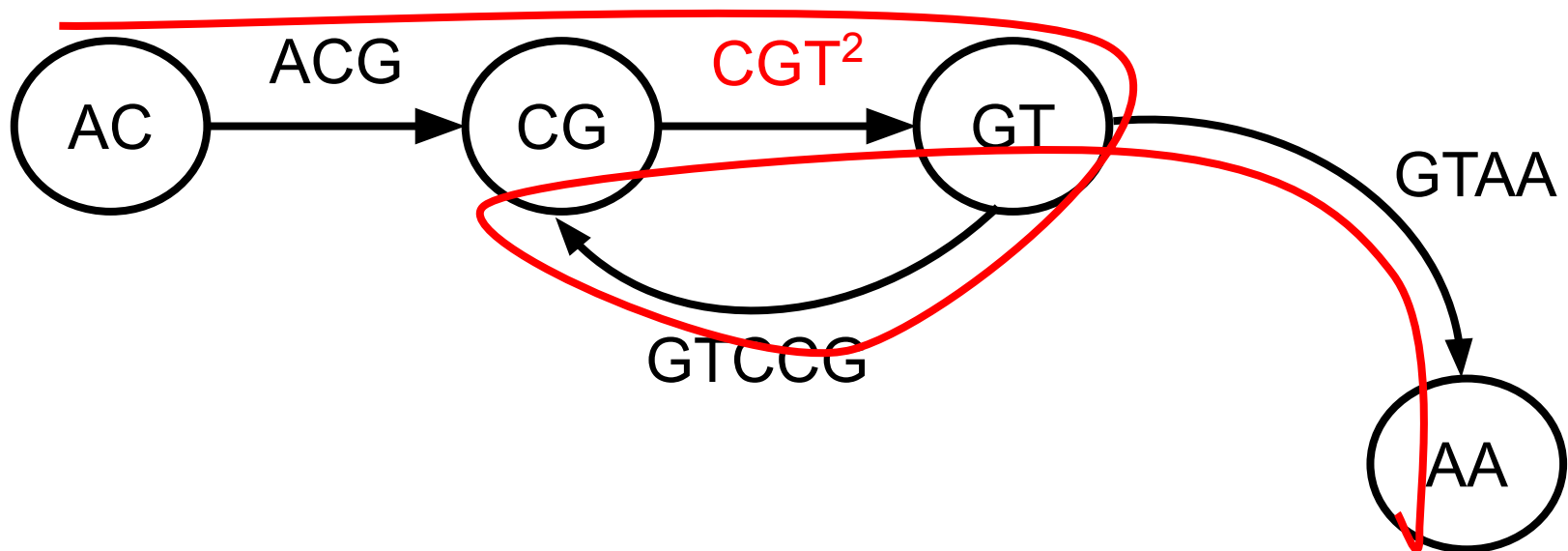k=2

# Eulerian path with multiplicities

**A<span style="color:red">CGT</span>C<span style="color:red">CGT</span>AA**

k=2

# More examples

**CCGTTG**

**TGCAGG**

**GTTGCA**

k=3

# More examples

CCG

**CCG**TTG
**TGCAGG**
**GTTGCA**

k=3

# More examples

CCG

CGT

**C<span style="color:red">CGT</span>TG**

**TGCAGG**

**GTTGCA**

k=3

# More examples

CCG   CGT

GTT

**CC<span style="color:red">GTT</span>G**
**TGCAGG**
**GTTGCA**

k=3

# More examples

CCG    CGT

GTT    TTG

**CCG<span style="color:red">TTG</span>**
**TGCAGG**
**GTTGCA**

k=3

# More examples

CCG    CGT

GTT    TTG

TGC    GCA

CAG    AGG

**CCGTTG**
**TGCAGG**
**GTTGCA**

k=3

# More examples

CCG

CGT

**CCGT**TG
**TGCAGG**
**GTTGCA**

k=3

CCGT

GTT

TTG

TGC

GCA

CAG

AGG

# More examples

CCG    CGT

**C<span style="color:red">CGTT</span>G**
**TGCAGG**
**GTTGCA**
k=3

CCGT

CGTT

GTT    TTG

TGC    GCA

CAG    AGG

68

# More examples

CCG    CGT

**CCGTTG**
**TGCAGG**
**GTTGCA**

k=3

CCGT

CGTT

GTT    TTG

GTTG

TGC    GCA

CAG    AGG

# What about real data?

CCG     CGT

GTT     TTG

**CCGTTG**
**TGCAGG**
**GTTGCA**

k=3

TGC     GCA

CAG     AGG

CCGT

CGTT

GTTG

TGCA

GCAG

CAGG

TTGC

# More examples

CCG → CGT

**CCGTTG**
**TGCAGG**
**GTTGCA**

k=3

GTT    TTG

TGC    GCA

CAG    AGG

CCGT

CGTT

GTTG

TGCA

GCAG

CAGG

TTGC

# What about real data?

CCG → CGT

CGT → GTT

GTT   TTG

**CCGTTG**
**TGCAGG**
**GTTGCA**

k=3

TGC   GCA

CAG   AGG

CCGT

CGTT

GTTG

TGCA

GCAG

CAGG

TTGC

**CCGTTG**
**TGCAGG**
**GTTGCA**

k=3

73

# More examples

CCG → CGT

GTT → TTG

TGC → GCA

CAG → AGG

**CCGTTG**
**TGCAGG**
**GTTGCA**

k=3

CCGT

CGTT

GTTG

TGCA

GCAG

CAGG

TTGC

# More examples



CCGTTG
TGCAGG
GTTGCA

k=3

CCGTTGCAGG

# Does k-mer size matter?

# One more example

TTTCATTC        AACGGGCA

AGCTTTTC        CTGCAACG

GGGCAATA        TGACTGCA

CATTCTGA

**K = 3**

# One more example
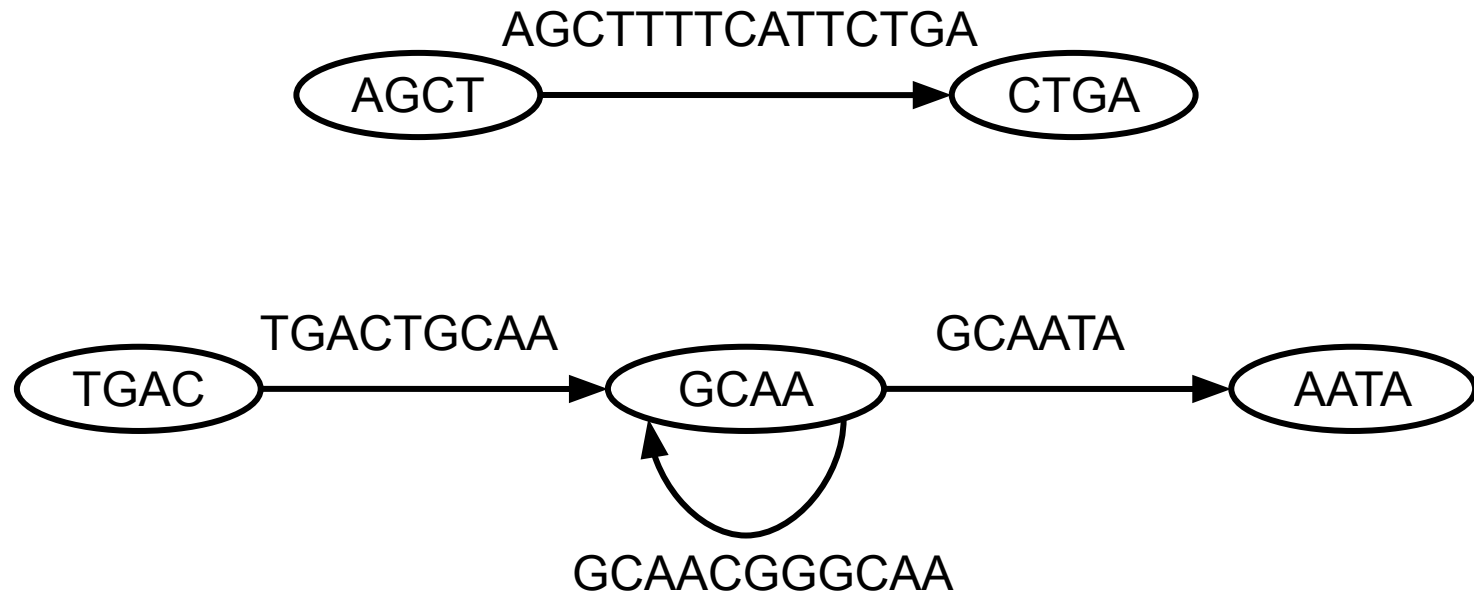
# One more example
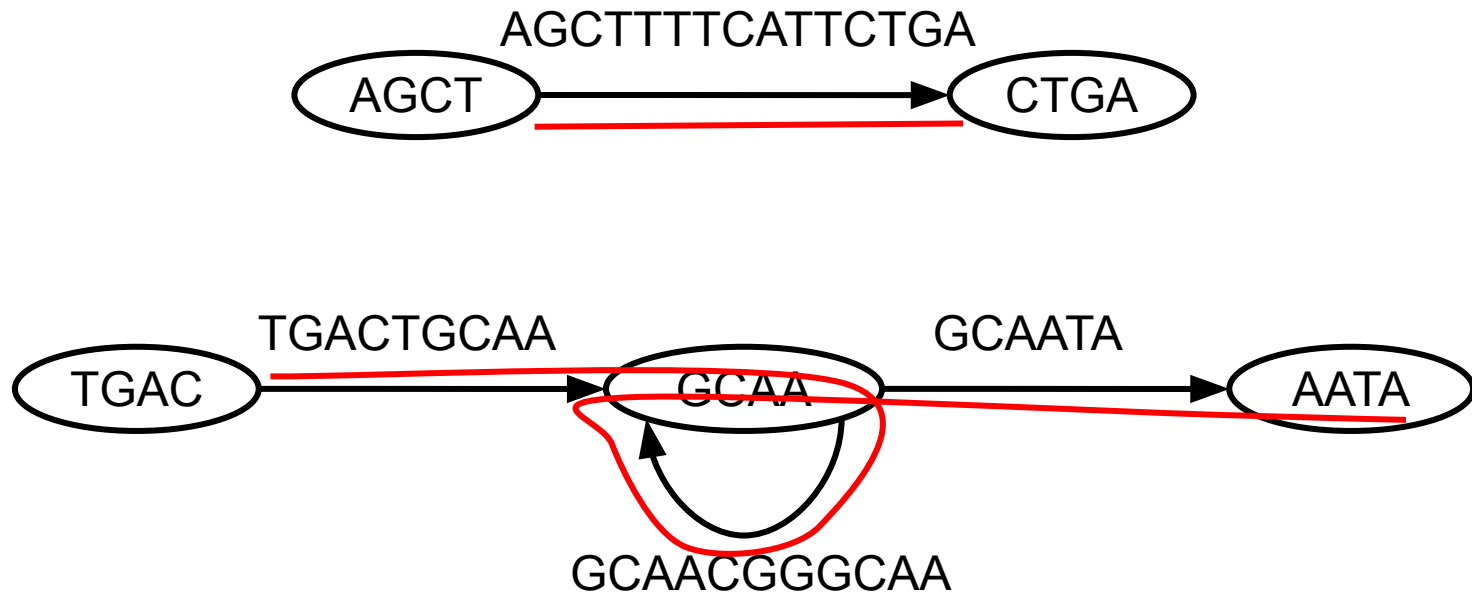
# One more example



**AGCTTTTCATTCTGACTGCAACGGGCAATA**

# One more example

TTTCATTC      AACGGGCA

AGCTTTTC      CTGCAACG

GGGCAATA      TGACTGCA
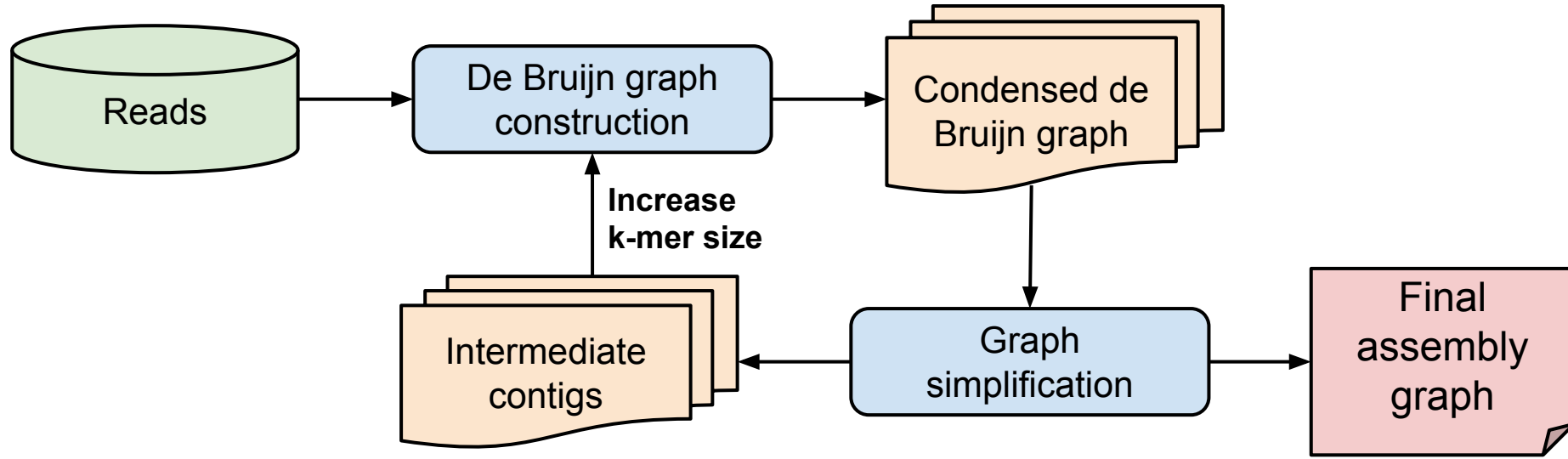
CATTCTGA

**K = 4**

# One more example

# One more example



**AGCTTTTCATTCTGACTGCAACGGGCAATA**

# How to select the *k*-mer size?

# How to select the *k*-mer size?

- Small k
  - Complex graph
  - Hard to resolve repeats
- Large k
  - Gaps in the assembly
- For normal data sets *k = ReadLength / 2 + ε*

# Iterative SPAdes run



- Smaller *k*-mer sizes are needed for reconstructing  low-coverage regions
- Larger *k*-mer sizes are needed for resolving short repeats

# DNA is double-stranded

# DNA is double-stranded

- Add k-mer and its reverse complement

- Use odd k to avoid self-complement vertices

  - rc(AA**T**TT) = AA**A**TT

  - rc(AATT) = AATT

# Removing sequencing errors

# Sequencing errors in de Bruijn graph

CCGTTG

CGTT**A**C

GTTGCA

TGCAGG

# Sequencing errors in de Bruijn graph

# Sequencing errors in de Bruijn graph

**CCGTTG**

**CGTTAC**

**GTTGCA**

**TGCAGG**

# How to remove a tip?

# How to remove a tip?

- Short length (usually less than 2 * k)
- Low coverage in respect to the main (correct path)
- Long length or high coverage — more likely to indicate a coverage gap

# More about sequencing errors

CCG → CGT → GTT

GTT → TTA → TAC → ACA → CAG → AGG

GTT → TTG → TGC → GCA → CAG

CCGTTG
CGTT**A**CAG
GTTGCA
TGCAGG

# More about sequencing errors



CCGTTG
CGTT**A**CAG
GTTGCA
TGCAGG

# More about sequencing errors

**CCGTTG**

**CGTT<span style="color:red">A</span>CAG**

**GTTGCA**

**TGCAGG**

# And what about bulges?

# And what about bulges?

- Erroneous path has lower coverage that correct one
- Rather small length
- In case of similar coverage or bigger — more likely to be result of diploidy

# Real life

# Real life

# Velvet assembler simplification

- Tip clipping

- Bulge removal

- Removing erroneous connections

# Velvet tip clipping

- Remove only if shorter than 2 * k

- Coverage is lower than of any alternative paths

- Iteratively process over the graph until no tips are left

# Velvet "tour bus" algorithm

- Distance between vertices A and B is

  $D(A, B) = length(E_{AB}) / coverage(E_{AB})$

  - Allows to go through reliable paths faster

- Start BFS from arbitrary node

- As soon as we came to already visited vertex

  - Align to alternative paths
  - Project low-covered path onto the main one

# Velvet "tour bus" algorithm

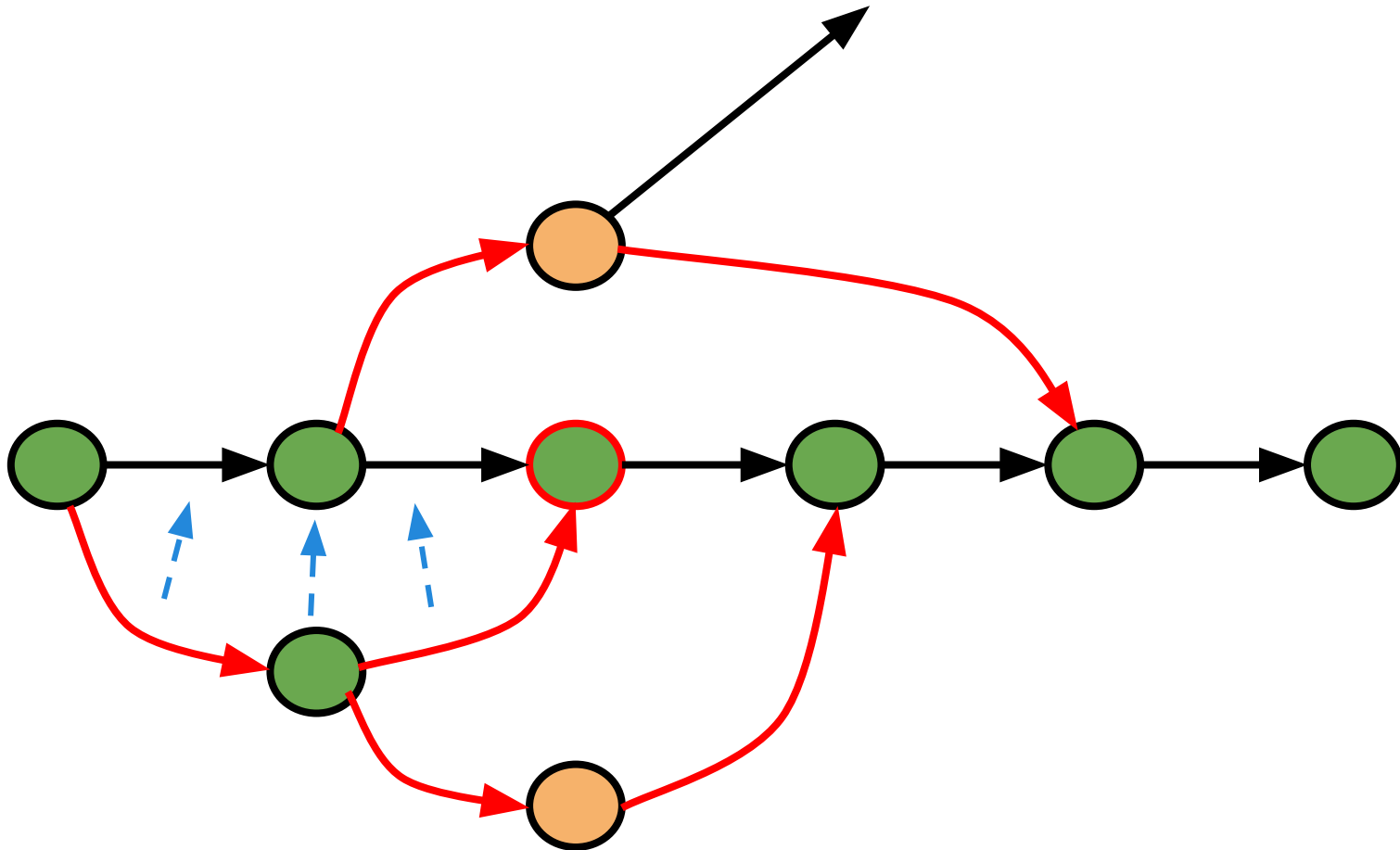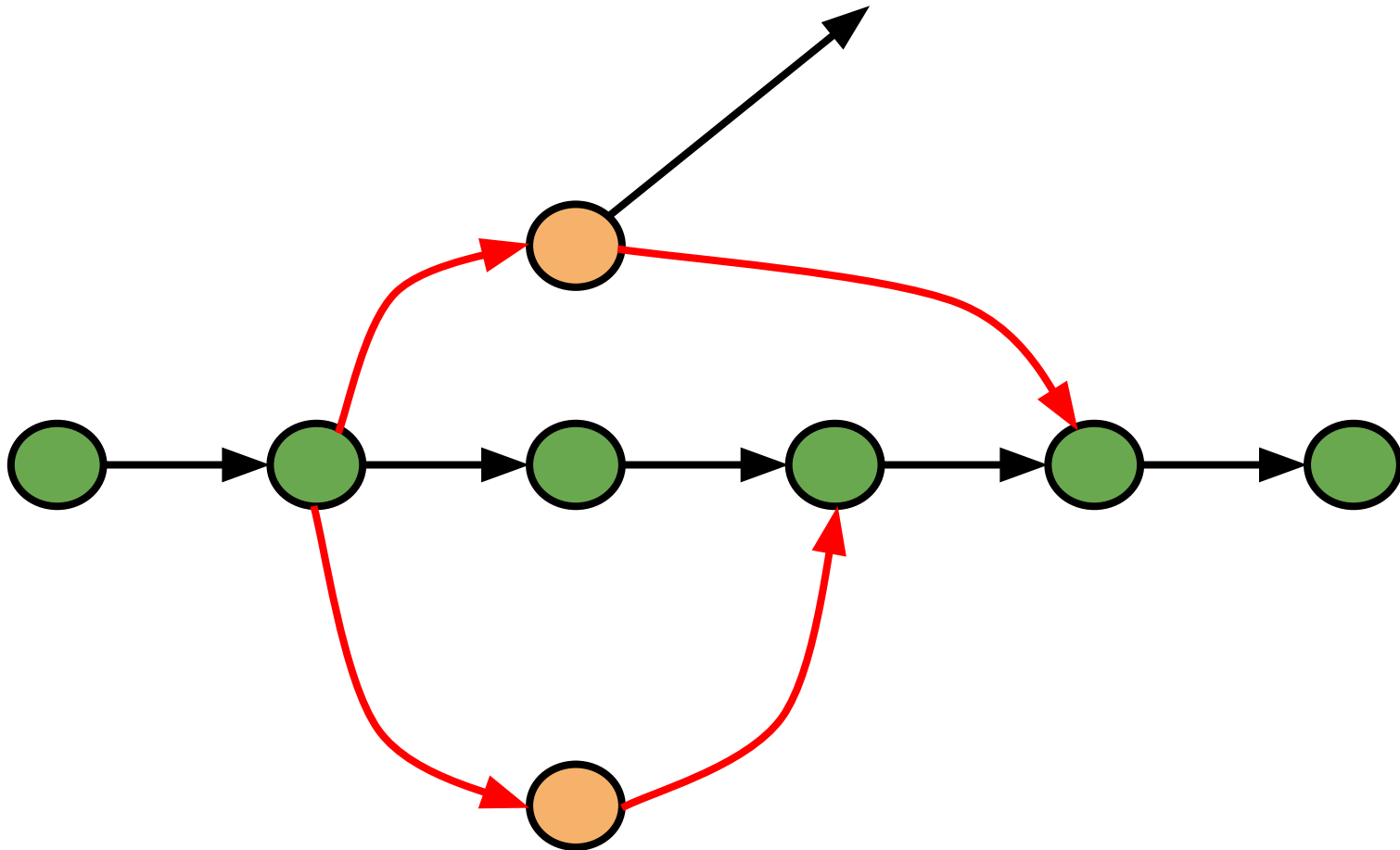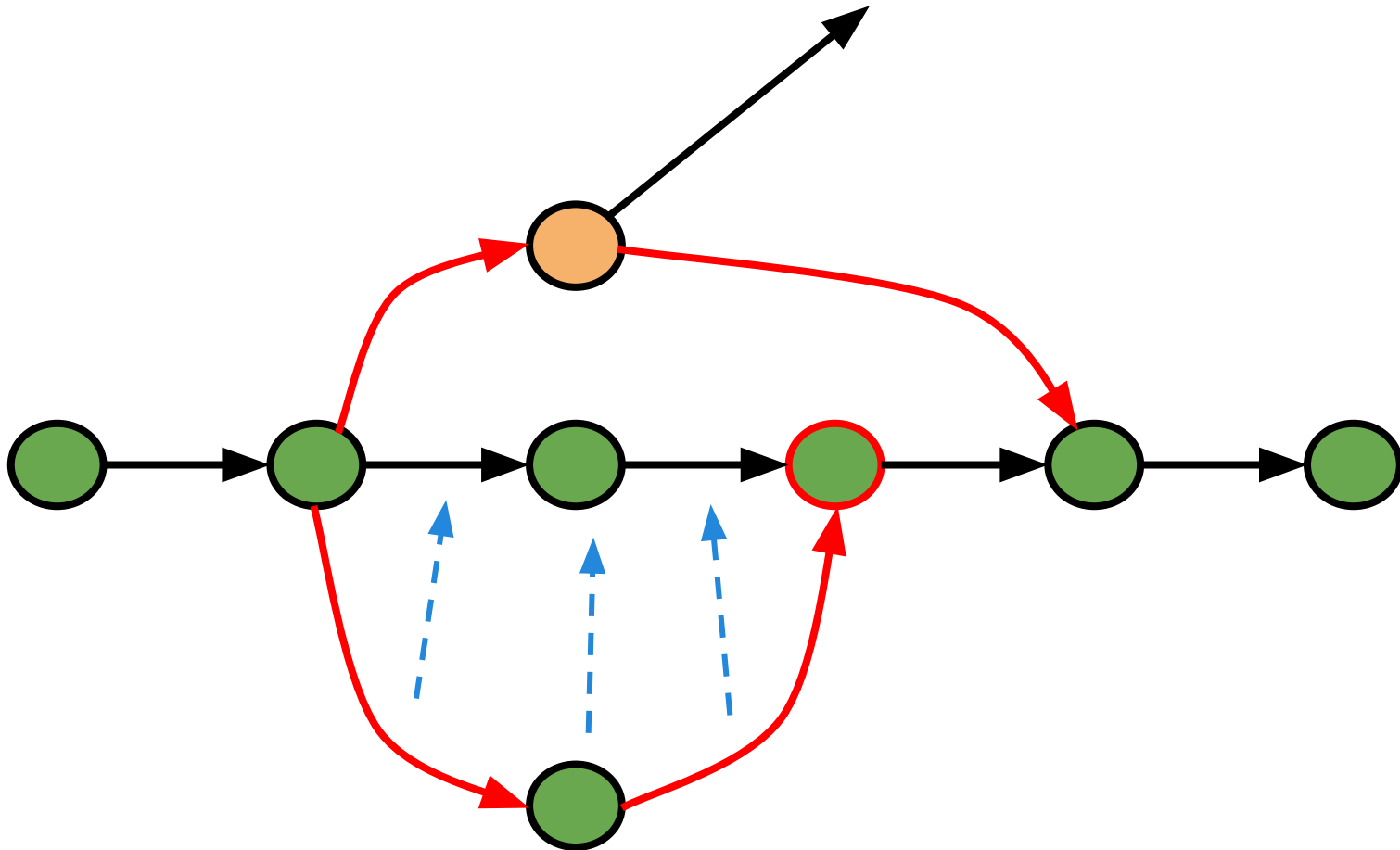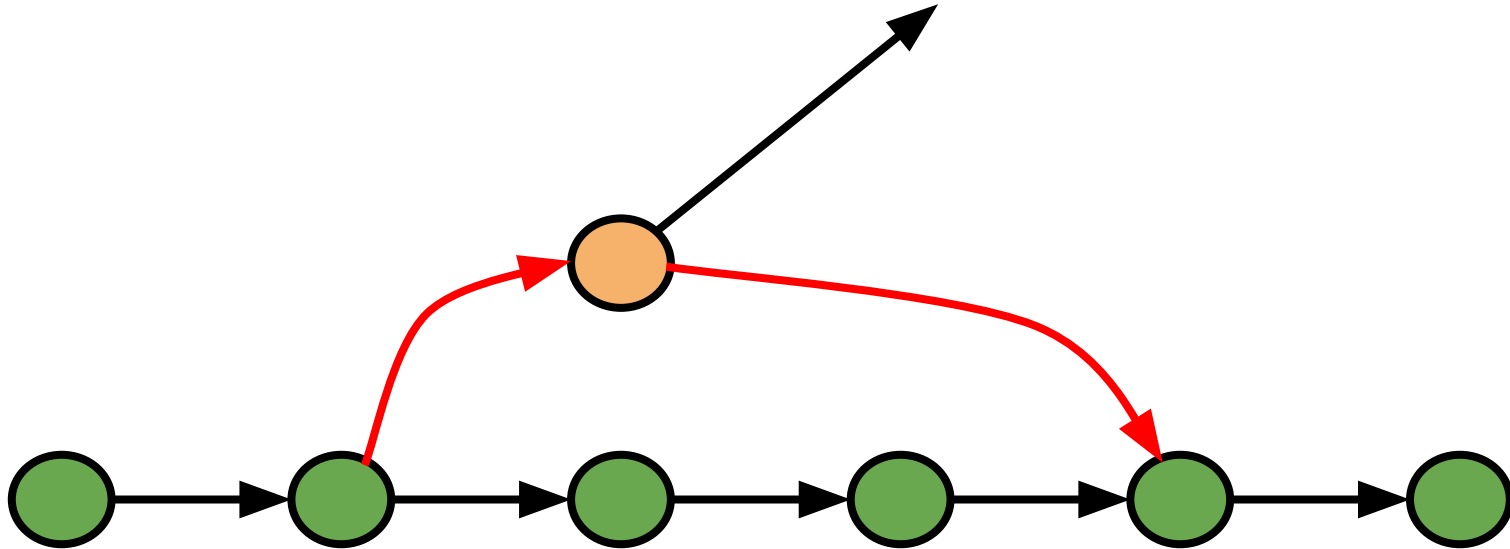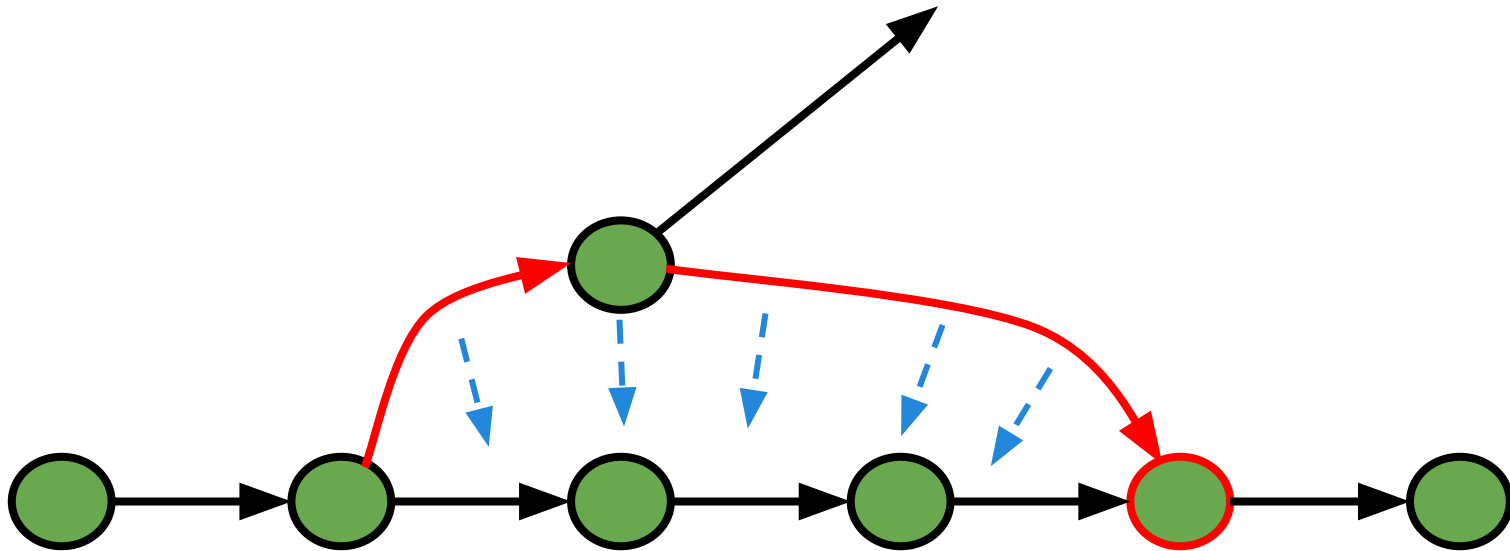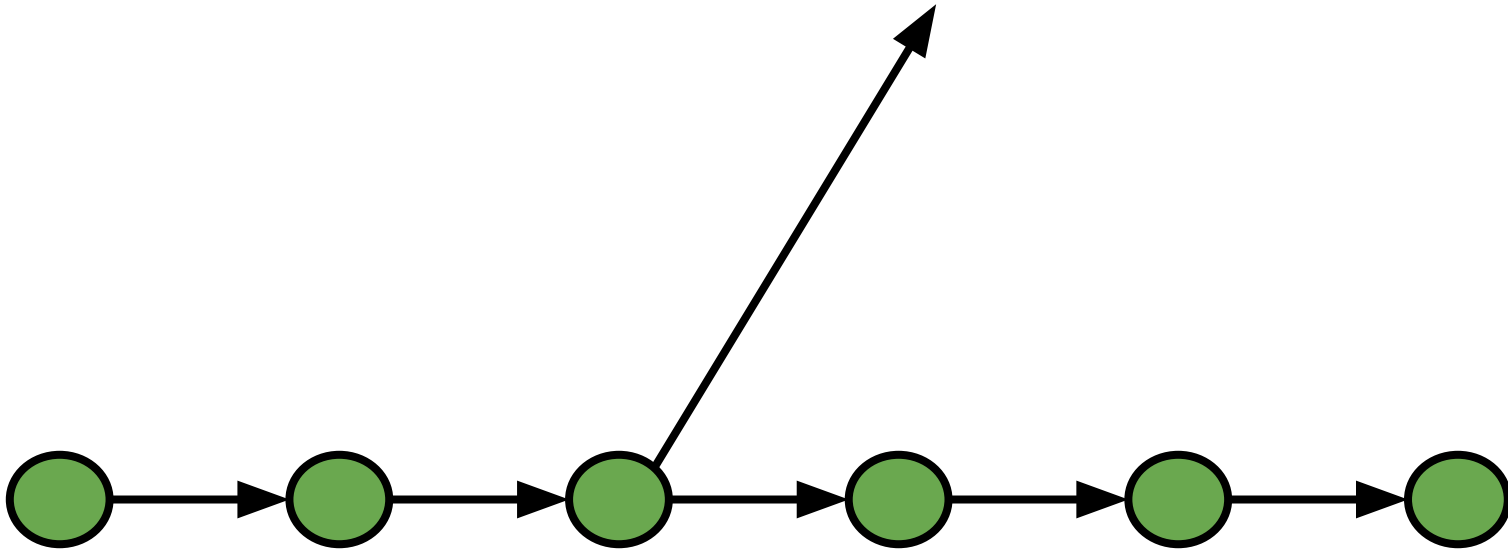# Velvet "tour bus" algorithm

# Velvet "tour bus" algorithm

# Velvet "tour bus" algorithm

# Velvet "tour bus" algorithm

# Velvet "tour bus" algorithm

# Velvet "tour bus" algorithm

# Velvet "tour bus" algorithm

# Velvet "tour bus" algorithm

# Velvet "tour bus" algorithm

# Velvet "tour bus" algorithm

# Velvet "tour bus" algorithm

# Velvet "tour bus" algorithm

# Velvet "tour bus" algorithm

# Velvet erroneous connection remover

- Erroneous connections don't have any recognized topological structure

- Have low coverage
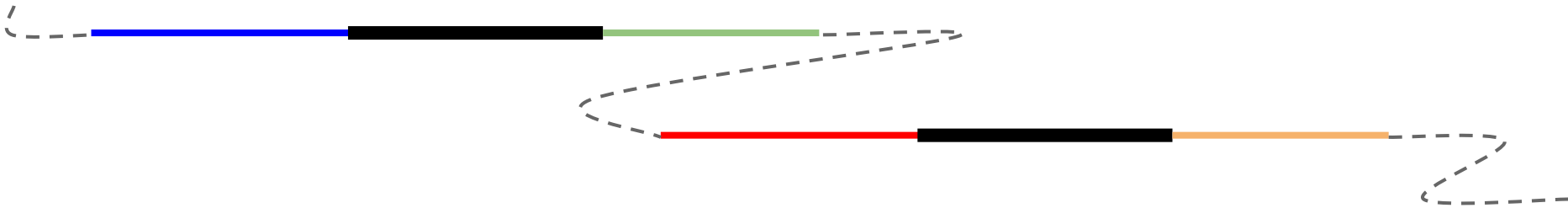
- Removed using simple coverage cutoff
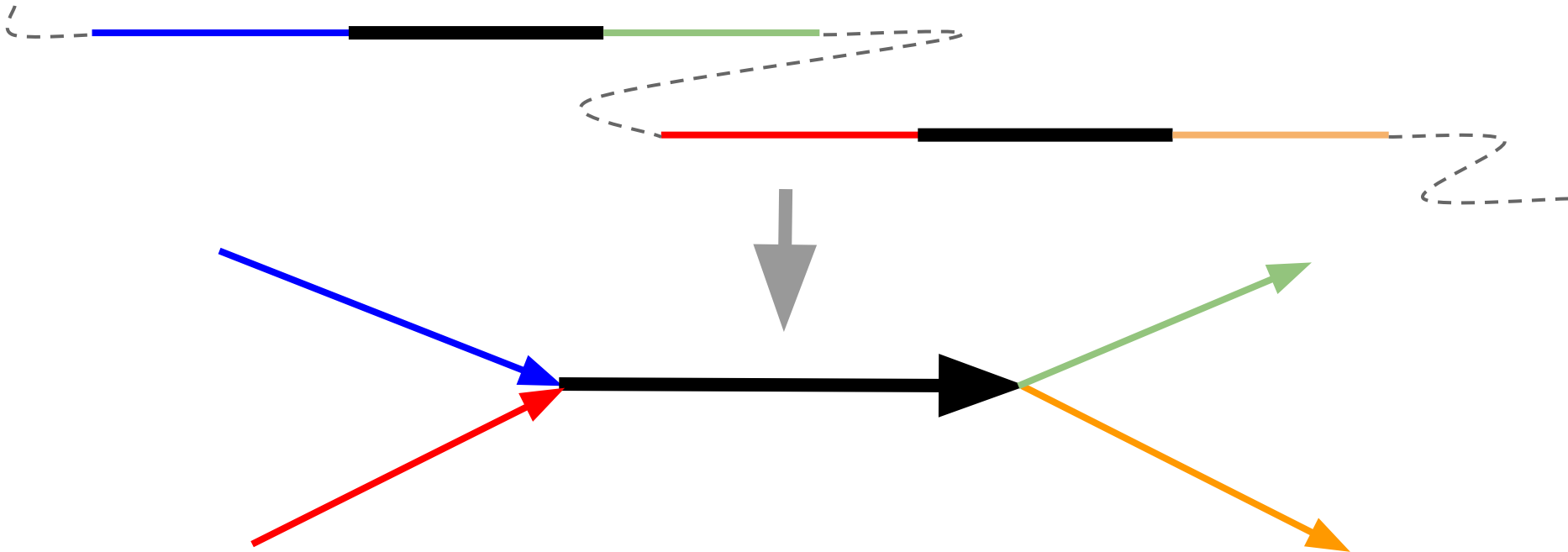
# Homework

Implement de Bruijn graph

- Construction from FASTA/FASTQ
- Condensation
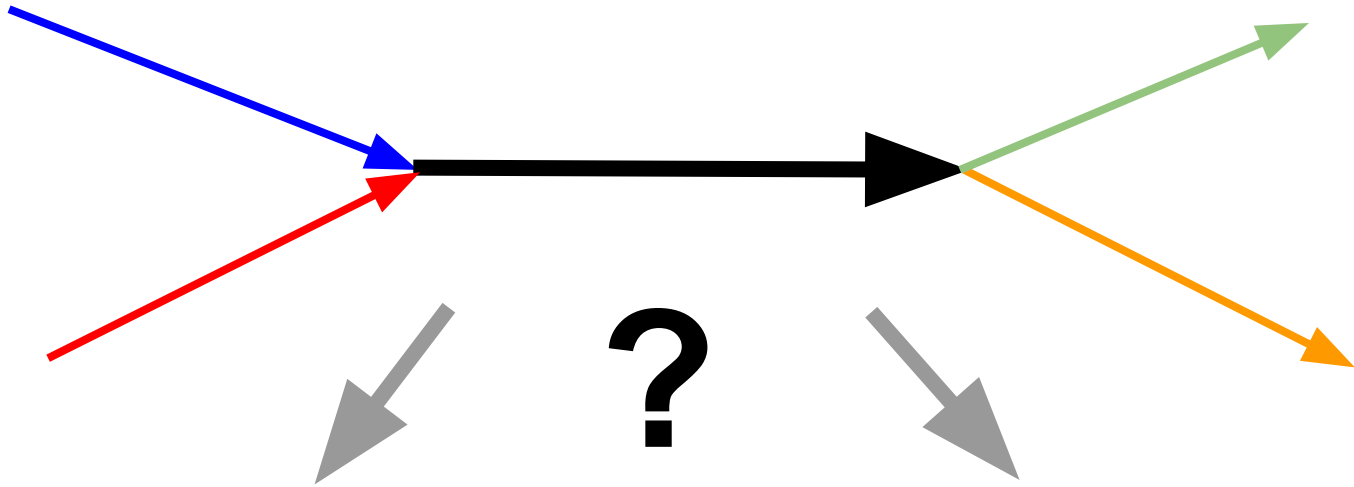- Output to FASTA/DOT/…
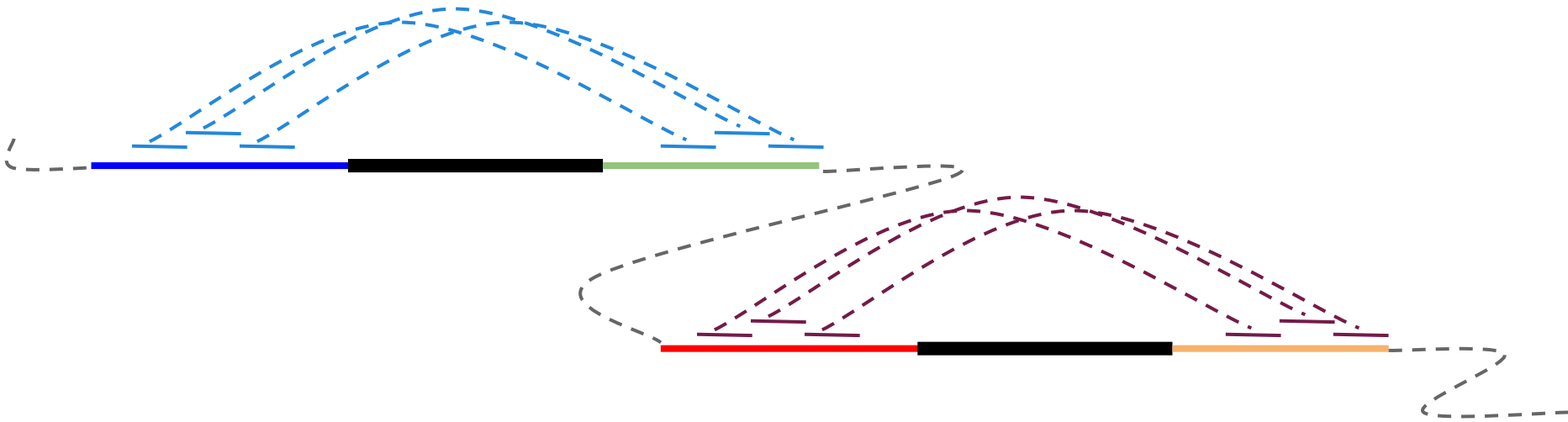- Simplification

# Resolving repeats
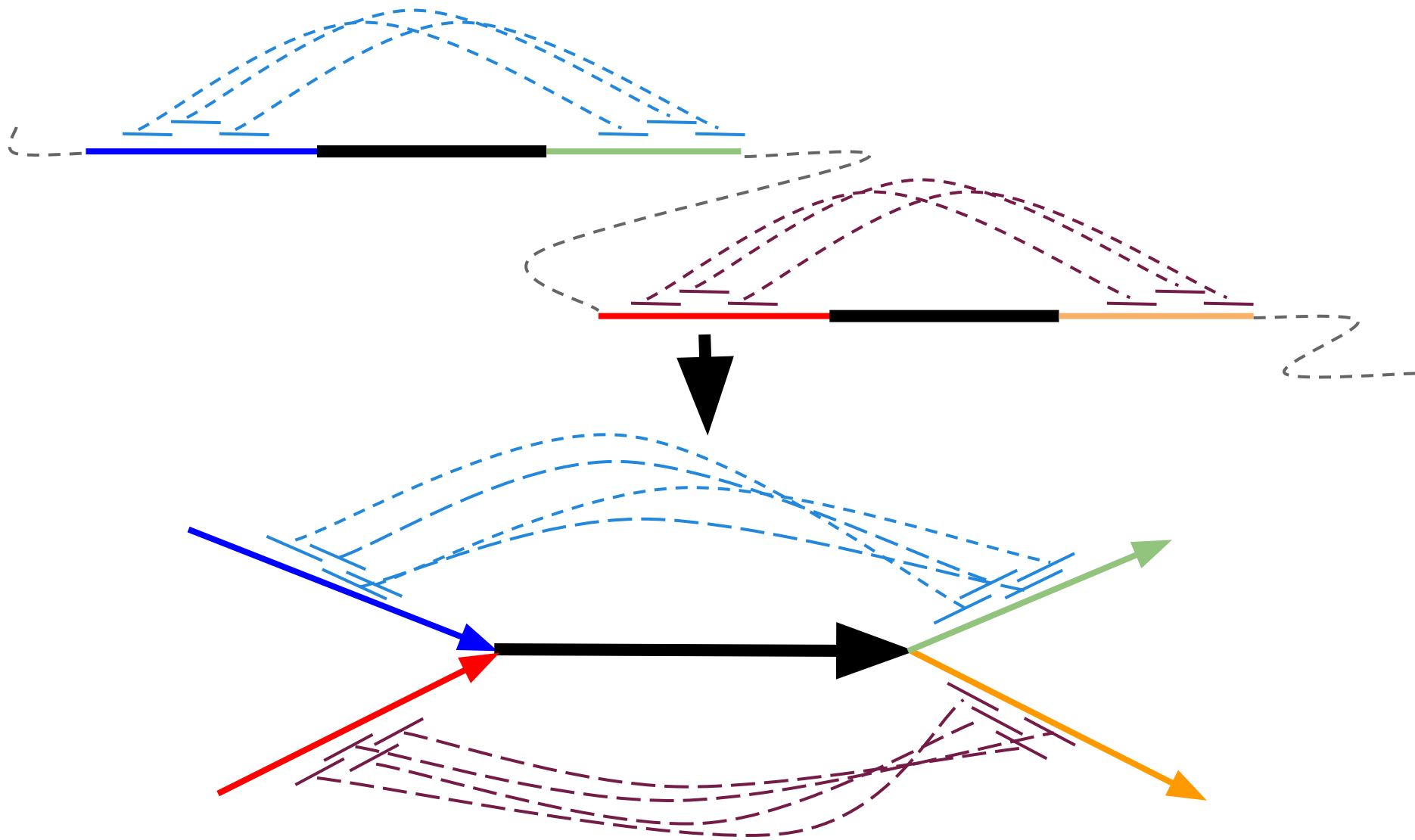
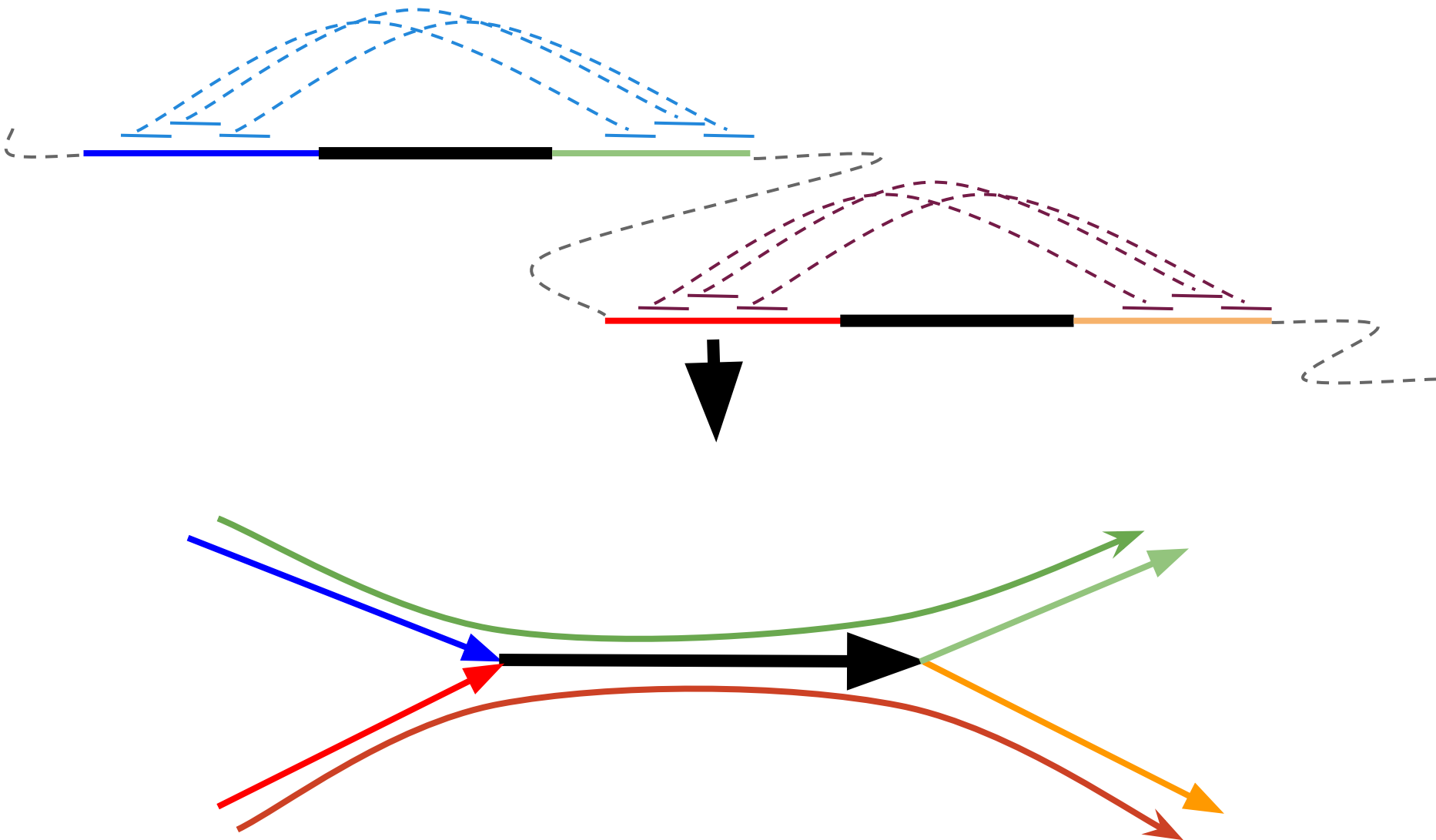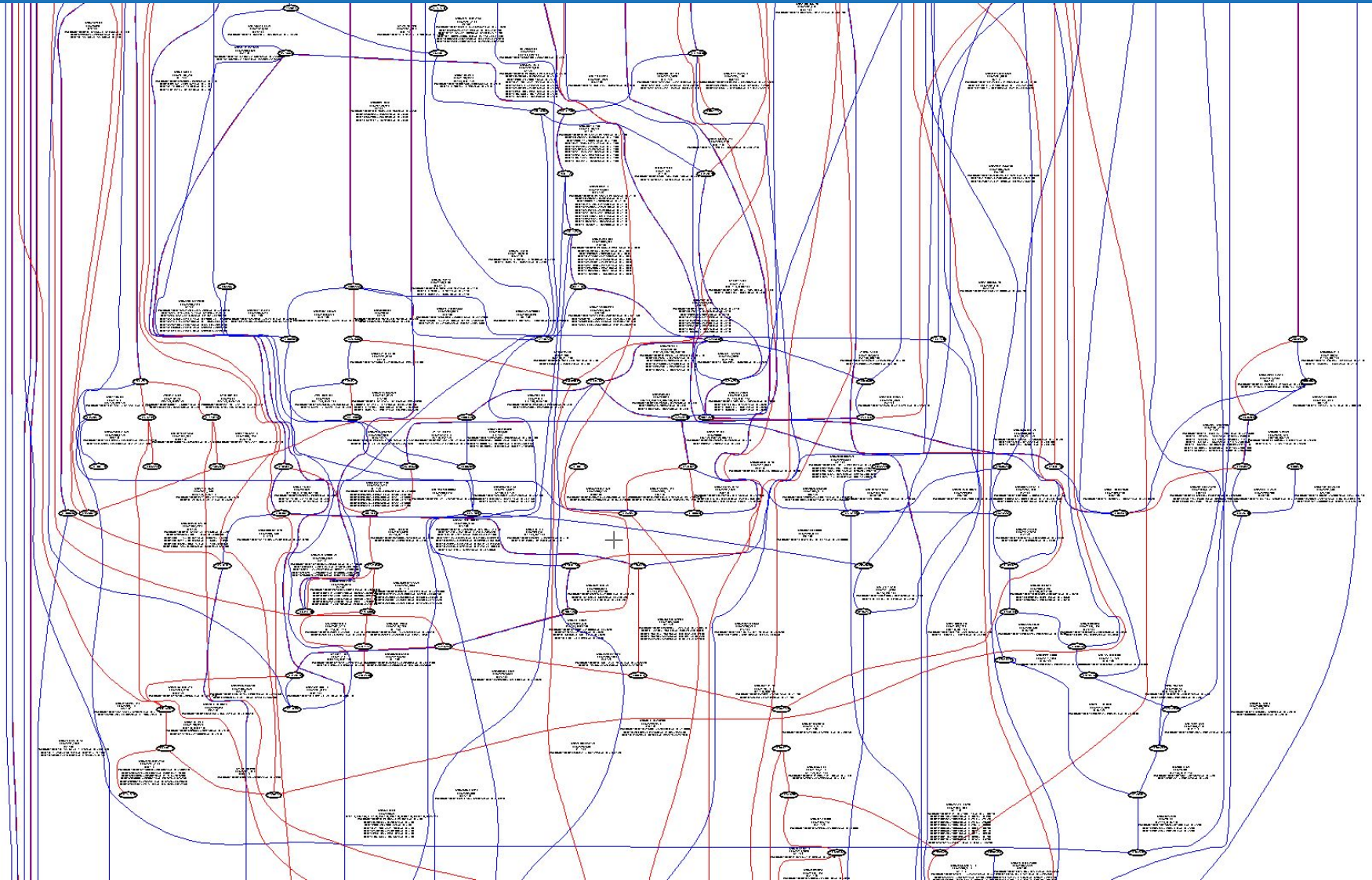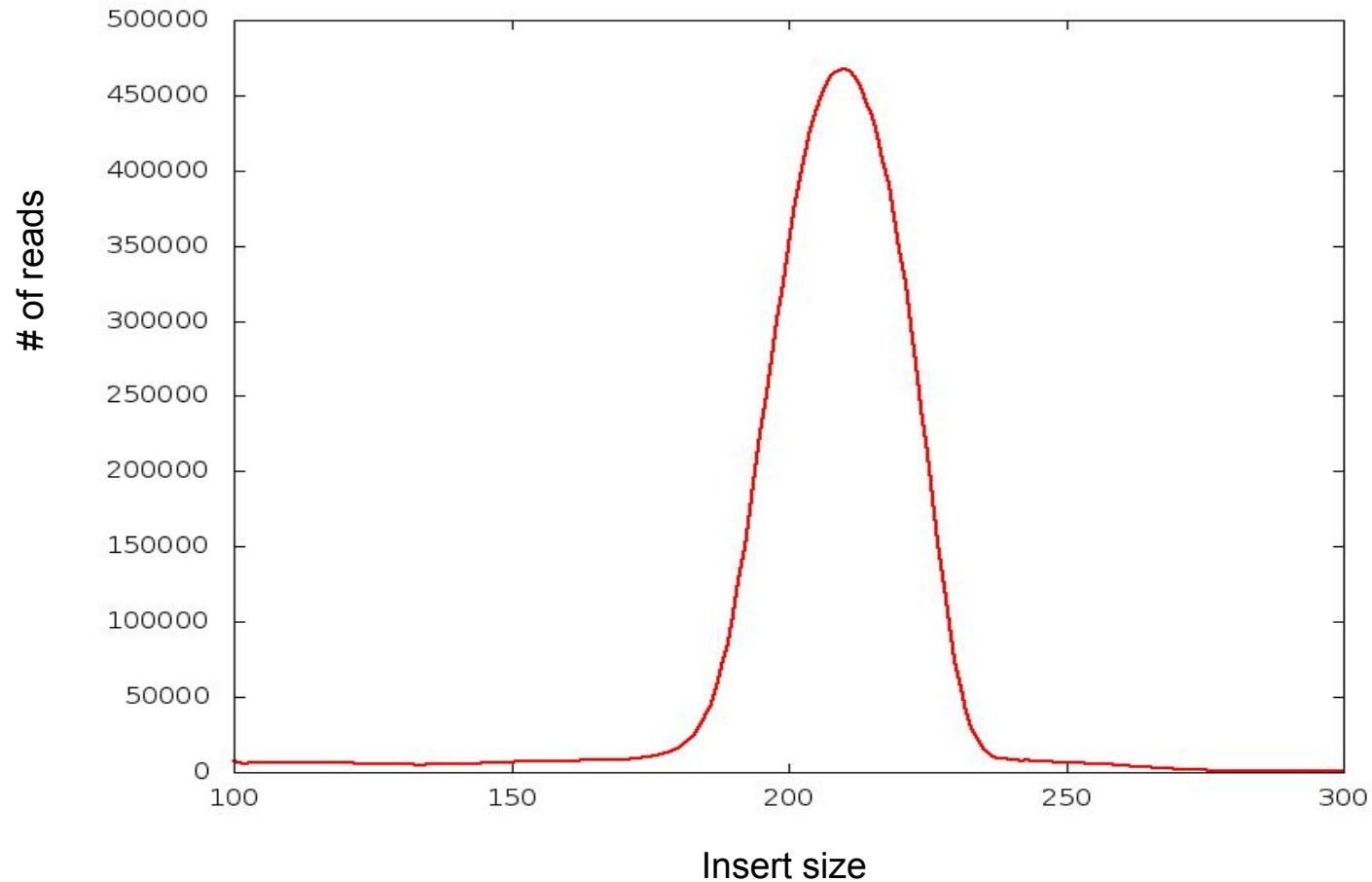# Resolving repeats

# Resolving repeats

# Resolving repeats

# Resolving repeats

# Resolving repeats

# Insert size distribution

## Paired-end reads

# Mate-pairs

B Outer End          Outer End B

3.5 kb DNA Molecule
with Biotin End Labels

# Mate-pairs



B Outer End — Outer End B

3.5 kb DNA Molecule
with Biotin End Labels

B Outer Ends

Circularized Molecule

# Mate-pairs



B Outer
  End

Outer B
End

3.5 kb DNA Molecule
with Biotin End Labels

B

Outer
Ends

Circularized Molecule

B

Outer
Ends

Internal
Fragment

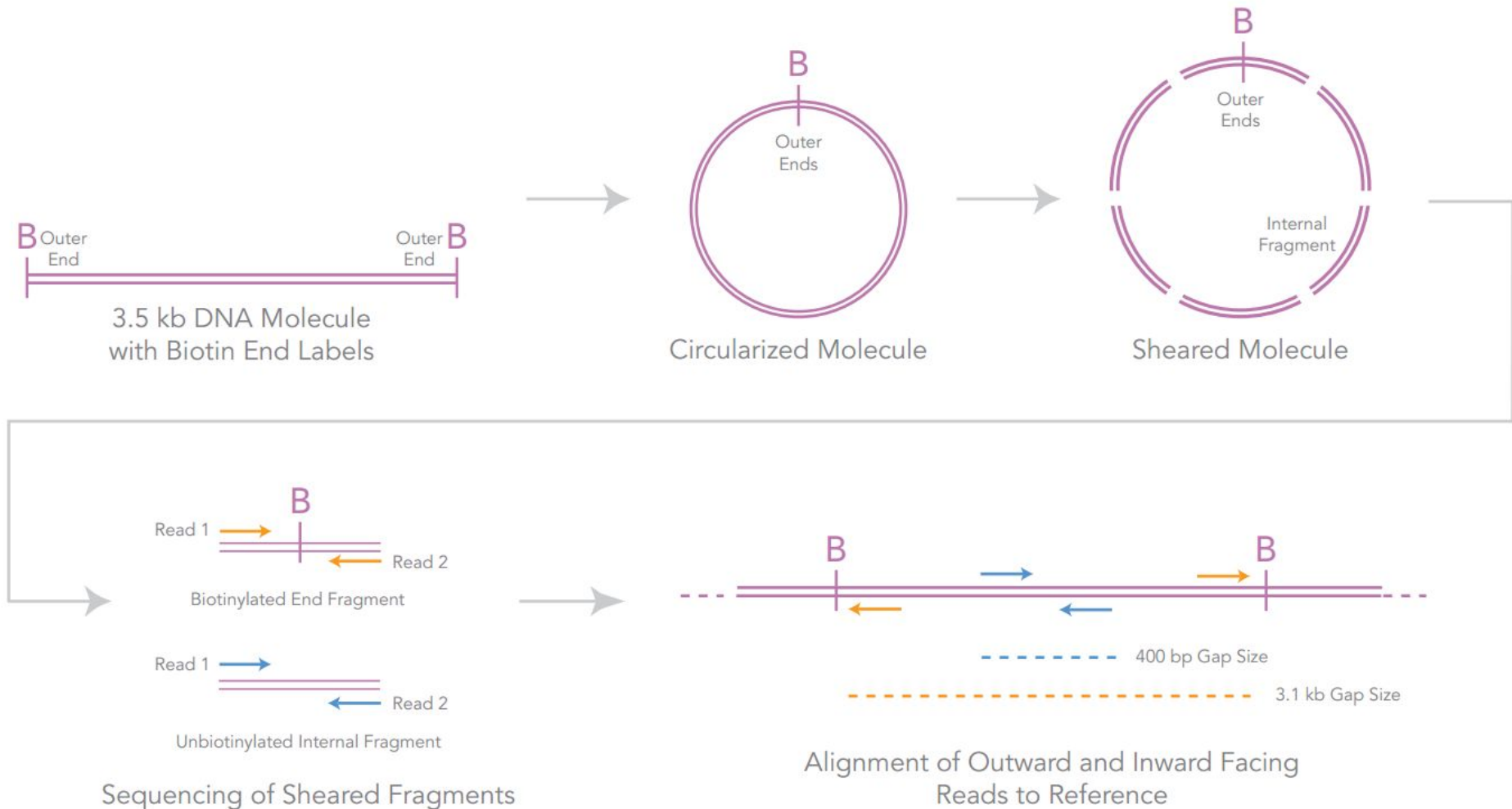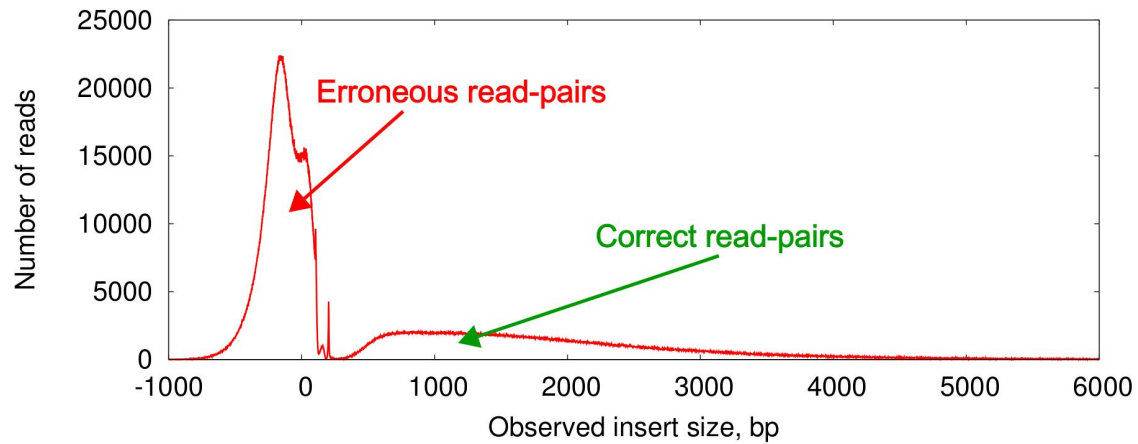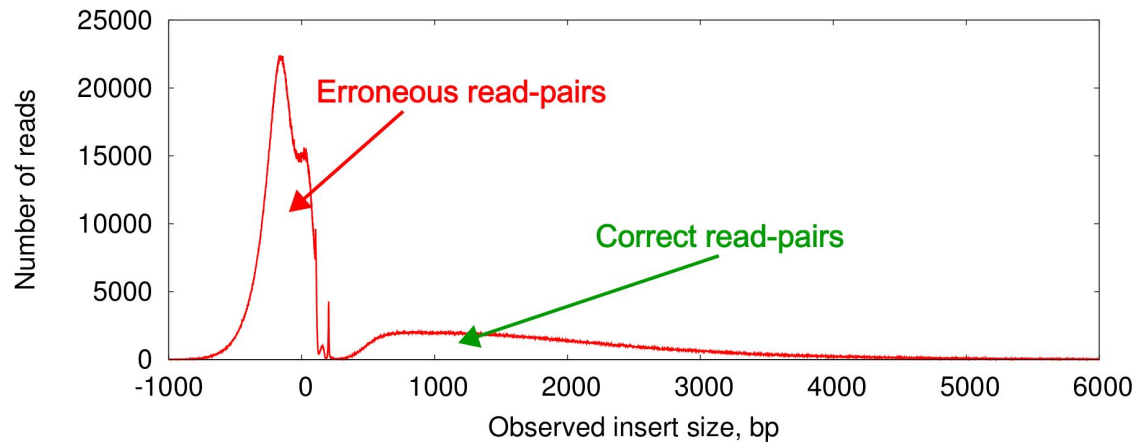Sheared Molecule

# Mate-pairs

# Mate-pairs

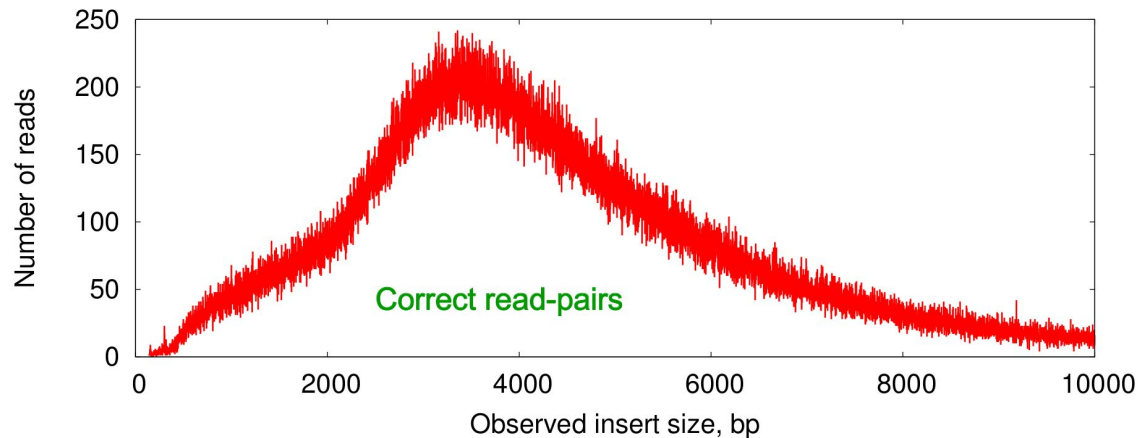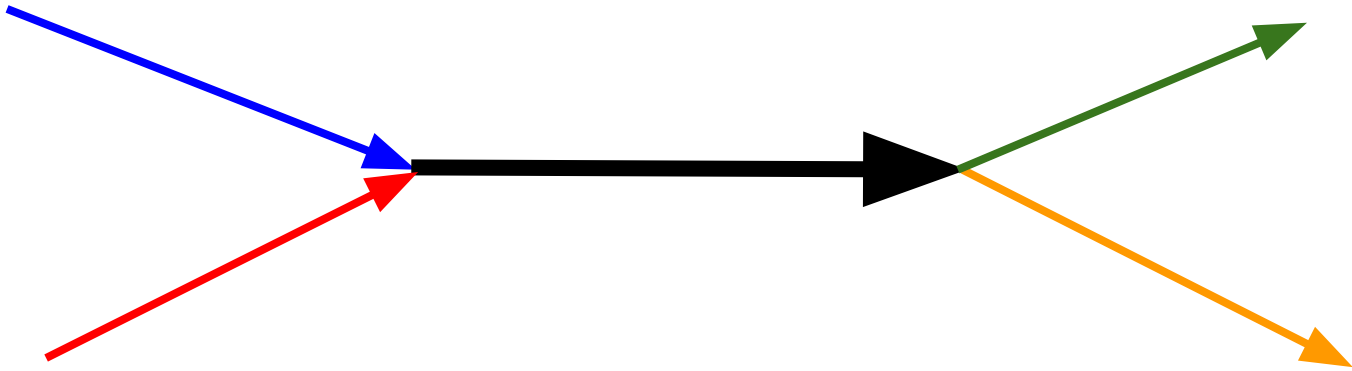# Mate-pairs

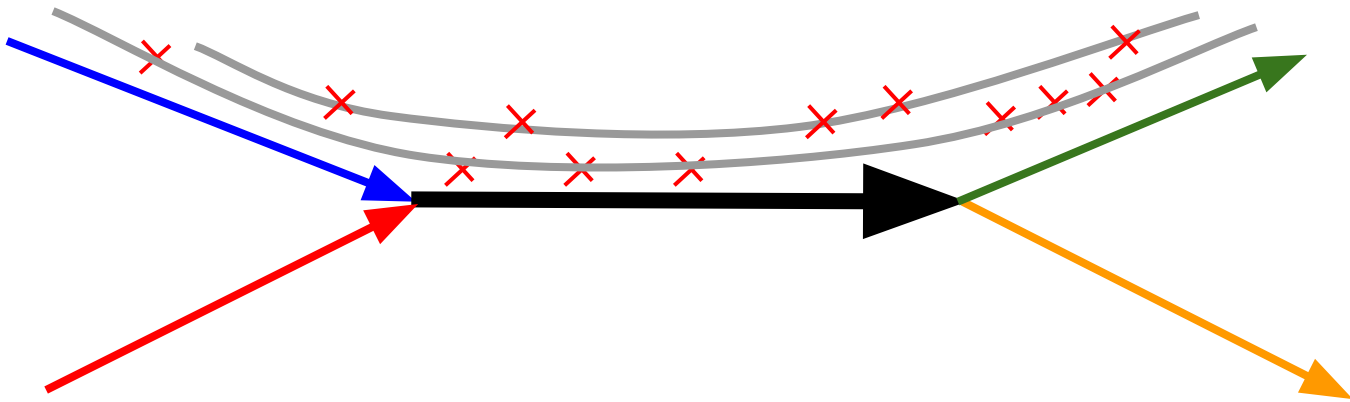Conventional mate-pairs:

# Mate-pairs
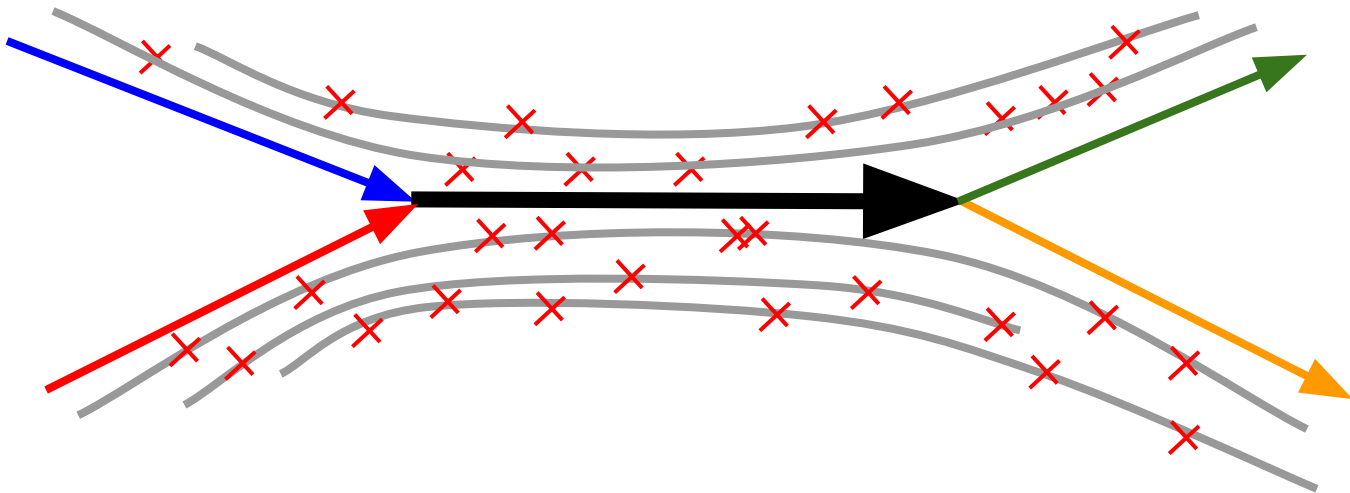
Conventional mate-pairs:
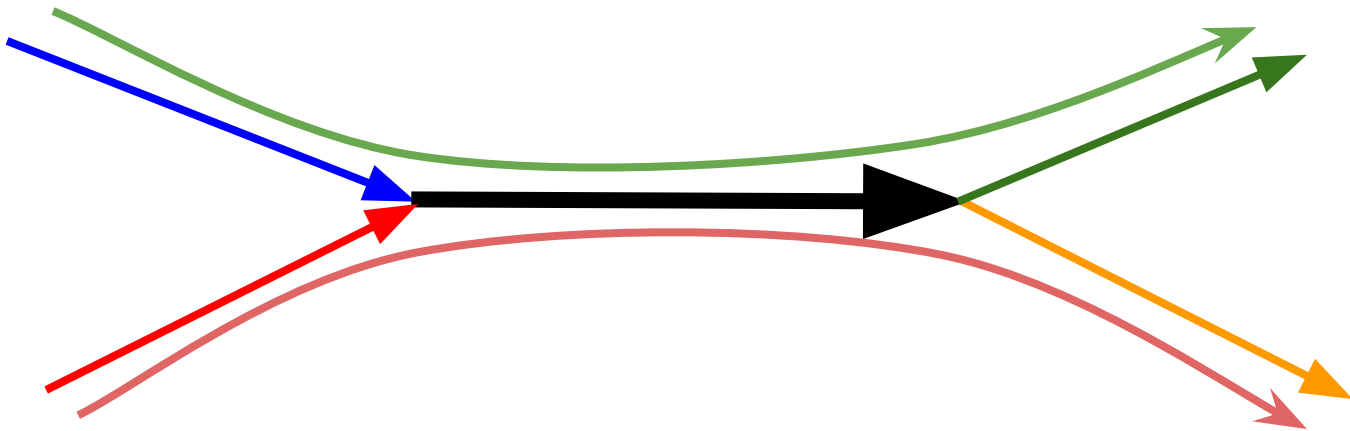


Illumina Nextera mate-pairs:

# Long reads to the rescue

# Long reads to the rescue

# Long reads to the rescue

# Sequencing technologies

| | SANGER SEQUENCING | illumina | Roche 454 SEQUENCING | ion torrent by life technologies | PACBIO | | NANOPORE |
|---|---|---|---|---|---|---|---|
| Protocols | | HiSeq, MiSeq | | | Subreads | CCS / HiFi | MinIon |
| Read length | 500-900 | 25-300 | 400-1100 | 200-400 | 20K-100K | 5K-20K | 1K-3M |
| Error rate | 0.001-0.1% | 0.1-1% | 1% | 1-2% | 2-10% | 0.1-2% | 5-15% |
| Error type | Indels & Mismatches | Mismatches only | Indels & Mismatches | Indels & Mismatches | Indels & Mismatches | Indels & Mismatches | Indels & Mismatches |
| Comments | Remains the golden standard | Error rate grows at the end of read | Problems with homopolymers | Problems with homopolymers | Errors distributed randomly | Error rate depends on sequencing settings | Typically several deletions in a row |
| Cost | $$$$$ | $$ | $$ | $$ | $$$ | $$$ | $ |

# PacBio only assembly

**Thm:**

**Perfect assembly possible iff**

**a) errors random**

**b) sampling is Poisson**

**c) reads long enough 2 solve repeats.**

**Note: e-rate not needed**

*Gene Meyers' twitter*

# New long reads vs Sanger assembly

High error rate => overlap detection is harder

- miniasm
- MHAP
- …

# Thank you!

**Questions?**