# REVIEW:

# ASSEMBLE A GENOME? GENERAL STRATEGIES

| Genome size | Unlimited $$ | Typical |
|---|---|---|
| >10Mb | | |
| 10Mb - 100Mb | | |
| > 100 Mb | | |

# ASSEMBLY

- OLC Assembly

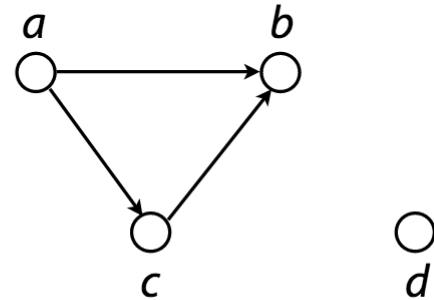| | |
|---|---|
| **Overlap** | **Build overlap graph** |
| Layout | Bundle stretches of the overlap graph into *contigs* |
| Consensus | Pick most likely nucleotide sequence for each contig |

# ASSEMBLY

Directed graph $G(V, E)$ consists of set of *vertices, V* and set of *directed edges, E*

Directed edge is an *ordered pair* of vertices.
First is the *source*, second is the *sink*.

    Vertex is drawn as a circle

    Edge is drawn as a line with an arrow
    connecting two circles



Vertex also called *node* or *point*

Edge also called *arc* or *line*

Directed graph also called *digraph*

$V = \{ a, b, c, d \}$

$E = \{ (a, b), (a, c), (c, b) \}$

Source    Sink

# ASSEMBLY – DE BRUIJN

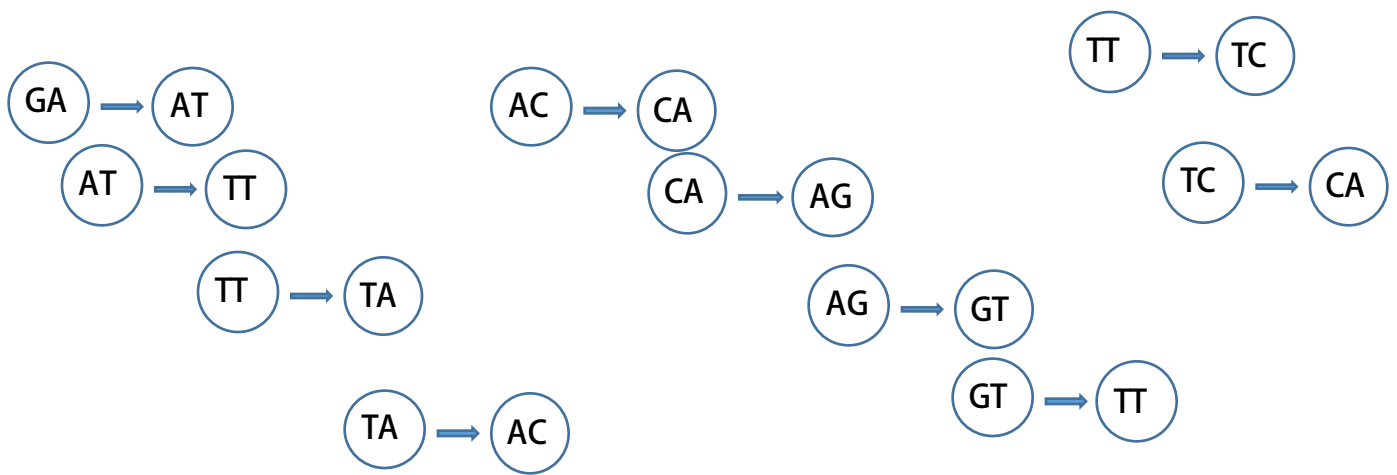Hamiltonian Path Problem

Eulerian Path Problem

# ASSEMBLY – DE BRUIJN

```
GATTAC
GAT
  ATT
   TTA
    TAC
```

```
ACAGTTCA
ACA
  CAG
   AGT
    GTT
     TTC
      TCA
```

# ASSEMBLY – DE BRUIJN

GAT   ATT   TTA   TAC   ACA   CAG   AGT   GTT  TTC  TCA

# ASSEMBLY – DE BRUIJN

GAT   ATT   TTA   TAC   ACA   CAG   AGT   GTT  TTC  TCA

GA → AT → TT → TA → AC → CA → AG → GT → TT → TC → CA

GAT   ATT   TTA   TAC   ACA   CAG   AGT   GTT  TTC  TCA
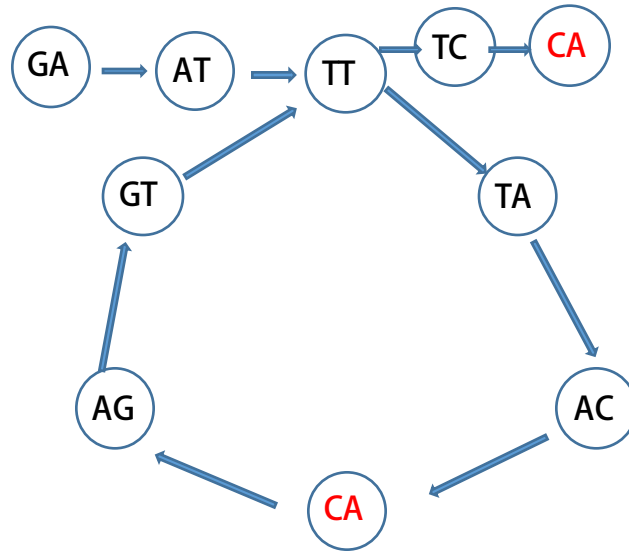
# ASSEMBLY – DE BRUIJN

GAT   ATT   TTA   TAC   ACA   CAG   AGT   GTT   TTC   TCA
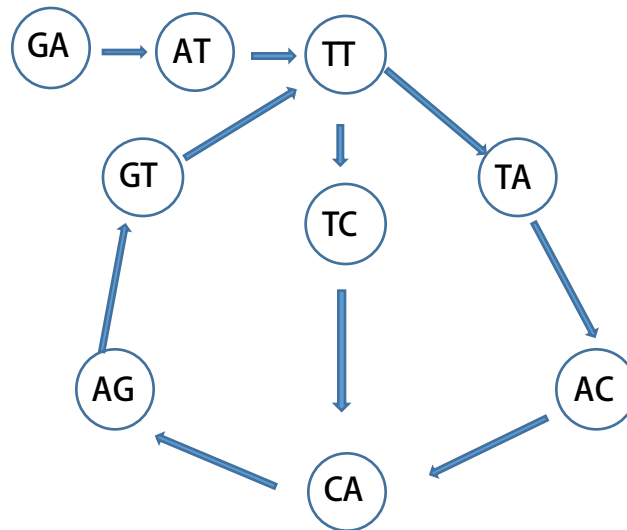
# ASSEMBLY – DE BRUIJN
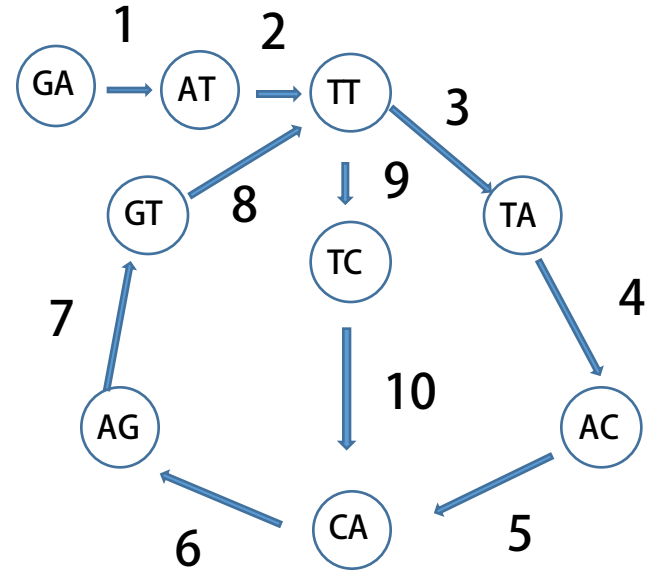
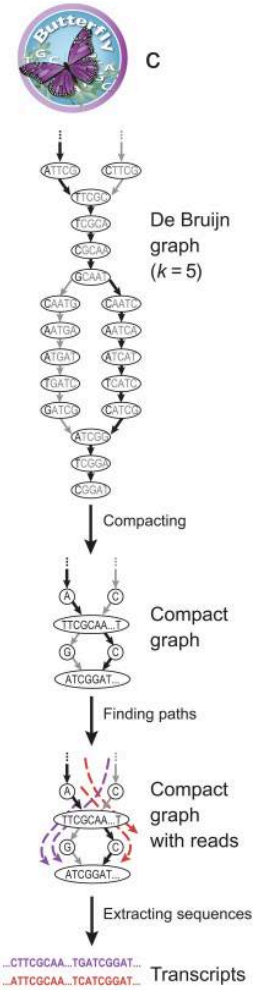GAT   ATT   TTA   TAC   ACA   CAG   AGT   GTT  TTC  TCA

# ASSEMBLY – DE BRUIJN

GAT  ATT  TTA  TAC  ACA  CAG  AGT  GTT TTC TCA
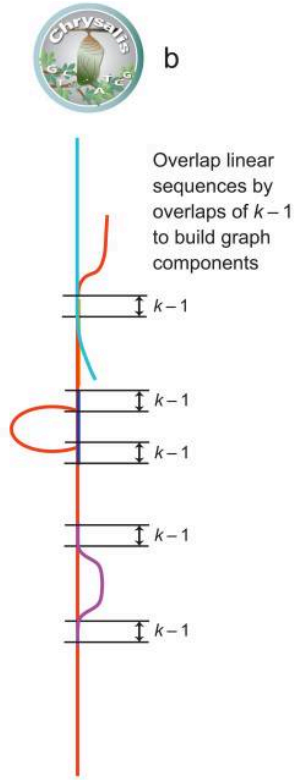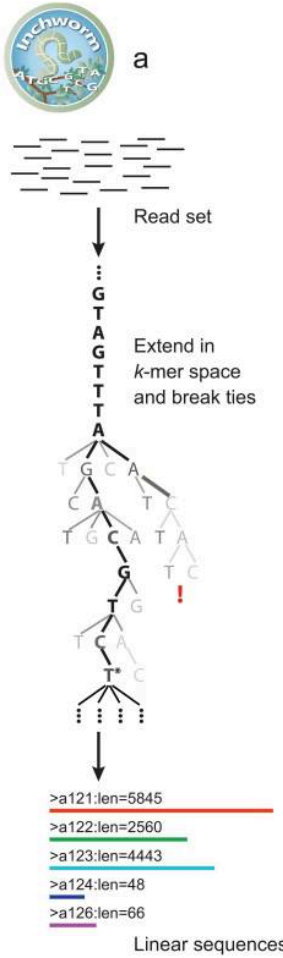
# TRANSCRIPTOME ASSEMBLY

Trinity



a

Read set

Extend in
k-mer space
and break ties

>a121:len=5845
>a122:len=2560
>a123:len=4443
>a124:len=48
>a126:len=66

Linear sequences

b

Overlap linear
sequences by
overlaps of k − 1
to build graph
components

$k − 1$

$k − 1$

$k − 1$

$k − 1$

$k − 1$

c

De Bruijn
graph
($k = 5$)

Compacting

Compact
graph

Finding paths

Compact
graph
with reads

Extracting sequences

...CTTCGCAA...TGATCGGAT...
...ATTCGCAA...TCATCGGAT...  Transcripts

# REVIEW:

Mapping

# MAPPING - BWT

|   | A | B | A | A | B | A |
|---|---|---|---|---|---|---|
| $ | a | b | a | a | b | a |
| a | $ | a | b | a | a | b |
| a | a | b | a | $ | a | b |
| a | b | a | $ | a | b | a |
| a | b | a | a | b | a | $ |
| b | a | $ | a | b | a | a |
| b | a | a | b | a | $ | a |

# MAPPING – SAM/BAM

# GENE EXPRESSION

Inter- versus intra-sample comparison

# GENE EXPRESSION

$$TPM_i = \frac{X_i}{\tilde{l}_i} * \left( \frac{1}{\sum_j \frac{X_j}{\tilde{l}_j}} \right) * 10^6$$