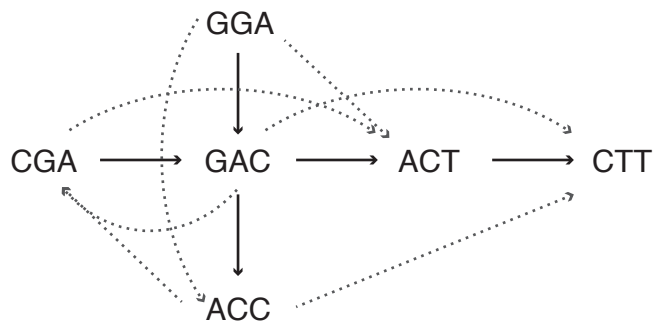


**a**  $K = \{CGA, GAC, ACT, CTT, GGA, ACC\}$

**b**



————→  $k-1$  long overlaps

.....→  $k-2$  long overlaps

de Bruijn graph:  
(vertex-centric, uni-directional)

————→

overlap graph:

————→ .....→ ...

**c**

| Name                    | Strings                      | Masked superstring   |   |
|-------------------------|------------------------------|--|---|
| 1. isolated $k$ -mers   | CGA, GAC, ACT, CTT, GGA, ACC | CGAGACACTCTTGGAAACC<br>100100100100100100                      | $\ell=18$<br>$r=6, o=6$                             |
| 2. unitigs              | CGA, GAC, GGA, ACC, ACTT     | CGAGACGGAACCACTT<br>1001001001001100                           | $\ell=16$<br>$r=5, o=6$                             |
| 3. simplitigs/SPSS      | CGACTT, GGA, ACC             | CGACTTGGAAACC<br>111100100100                                  | $\ell=12$<br>$r=3, o=6$                             |
| 4. matchtigs            | CGACTT, GGACC                | CGACTTGGACC<br>111100111100                                    | $\ell=11$<br>$r=2, o=7$                             |
| 5. shortest superstring | GGACCGACTT<br>               | GGACCGACTT<br>1110111100 ①<br>or 1110101100<br>or 1010111100 ② | $\ell=10$<br>$r=2, o=7$<br>$r=3, o=6$<br>$r=3, o=6$ |

**d**

|                                    |  |
|------------------------------------|--|
| enc0: "CGACTT\nGG<br>A\nACC"       | simplitigs/matchtigs delimited by EOL<br>(state-of-the-art)      |
| enc1: "GGACCGACTT"<br>"1110111100" | superstring + mask strings                                       |
| enc2: "GGAcCGACTt"                 | mask-cased superstring   |
| enc3: "GGACCGACTT"<br>[ 3, 1, 4 ]  | superstring + mask RLE   |
| enc4: "GGACCGACTT"<br>[ 3 ]        | superstring + mask sparse enc.<br>(positions of non-final zeros) |

**e**

①  $K = \{CGA, 2\times GAC, ACT, CTT, GGA, ACC\}$   
 $X = \{CCG\}$

②  $K = \{CGA, 1\times GAC, ACT, CTT, GGA, ACC\}$   
 $X = \{CCG, 1\times GAC\}$