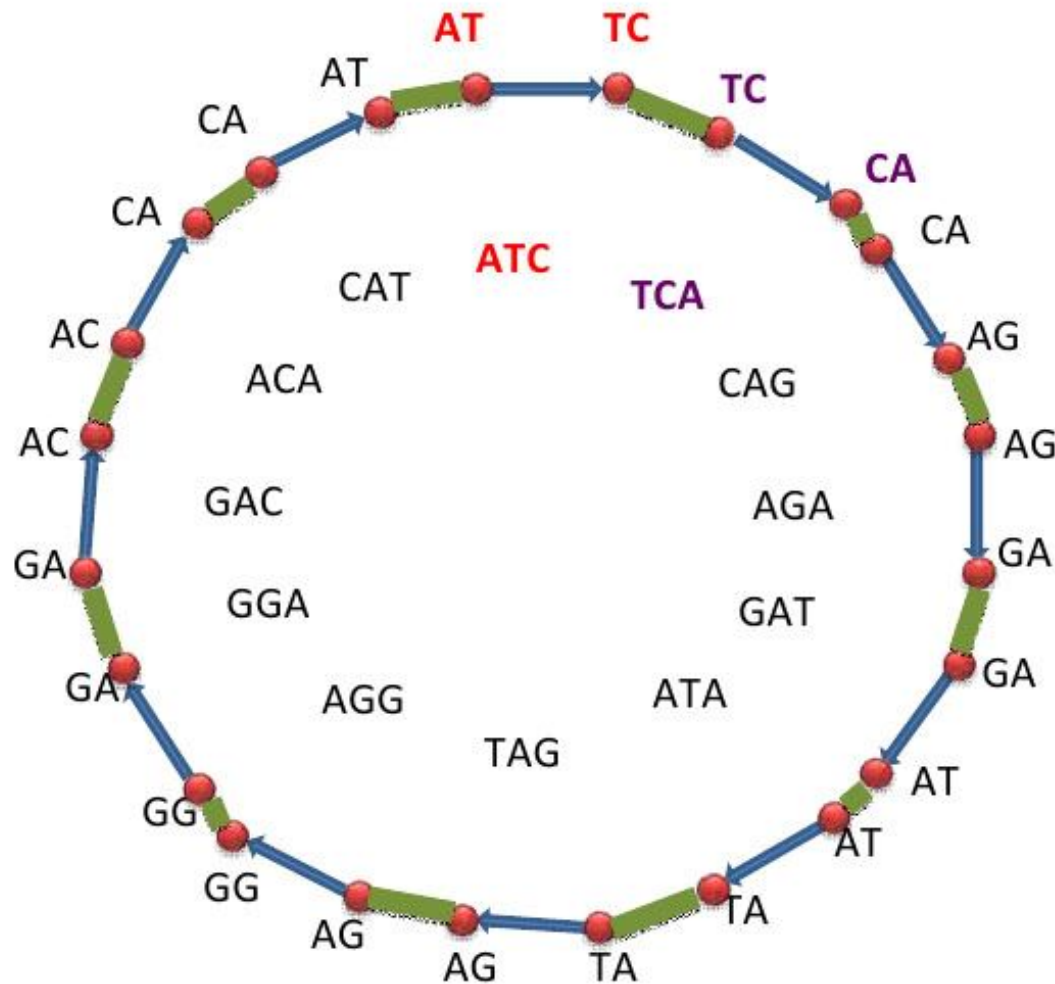


# Multiple libraries pathsets

Student: Nadiya Sitdykova  
Academic Adviser: Son Pham

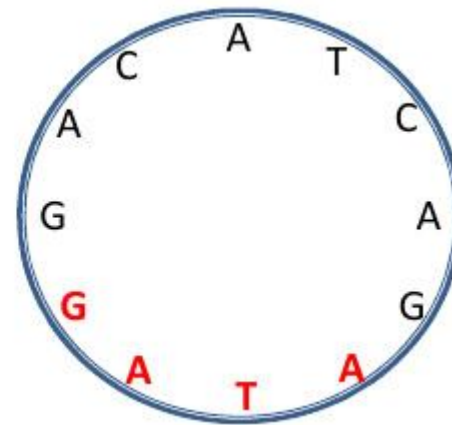
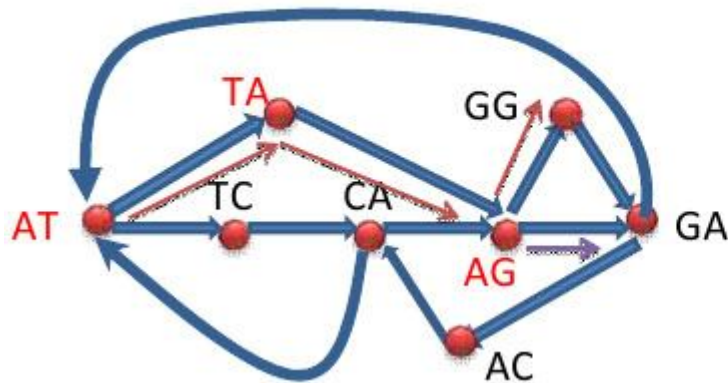
# Problem statement



- Genome is circular string.
- Reads cover all k-mers in the genome.
- Two libraries of read-pairs with different insert sizes  $d_1$  and  $d_2$ .

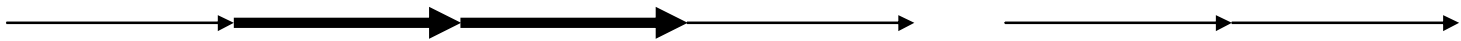
# Problem statement

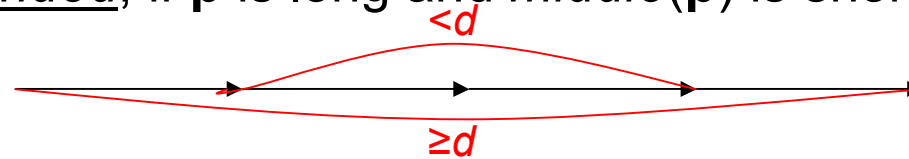
- Find the Eulerian cycle in the de Bruijn graph, corresponded to Genome from fundamental way (without heuristics).



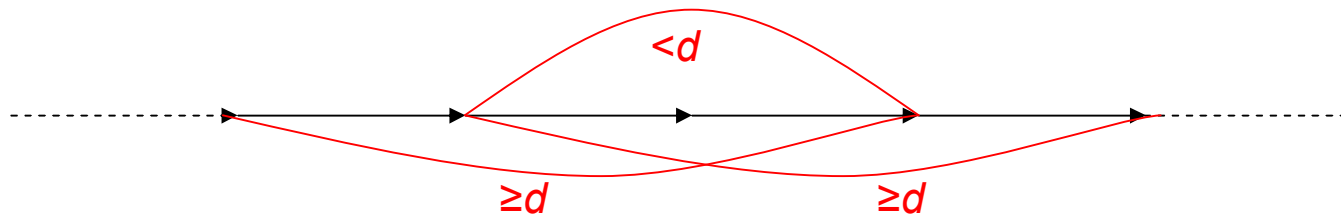
# Pathset assembly

## Definitions

- Path  $\mathbf{p}$  is long, if  $\text{length}(\mathbf{p}) \geq d$
- Path  $\mathbf{p}$  is short, if  $\text{length}(\mathbf{p}) < d$
- middle( $\mathbf{p}$ )  A horizontal line with arrows pointing right. The middle portion of the line is significantly thicker than the rest.
- Path  $\mathbf{p}$  is  $d$ -bounded, if  $\mathbf{p}$  is long and middle( $\mathbf{p}$ ) is short



- Path  $\mathbf{p}$  is  $d$ -maxpath, if  $\mathbf{p}$  is maximal  $d$ -bounded



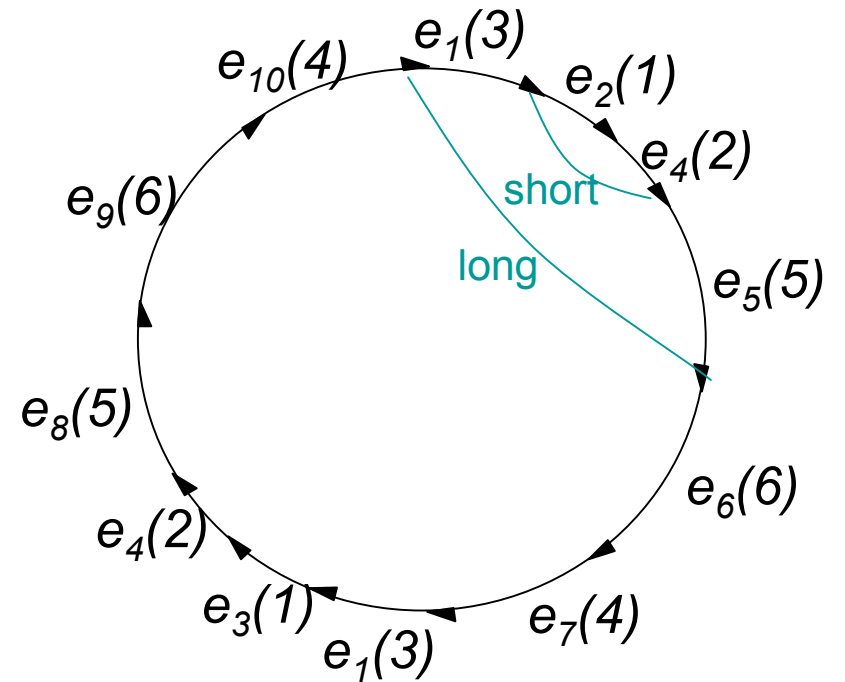
- First and last edges of  $d$ -bounded paths produce edge-pairs in the condensed de Bruijn graph.
- Biedge( $a|b,d$ ) defines Pathset( $a|b,d$ ) as all paths between  $a$  and  $b$  of length  $d$  in the condensed de Bruijn graph.

# Pathset assembly

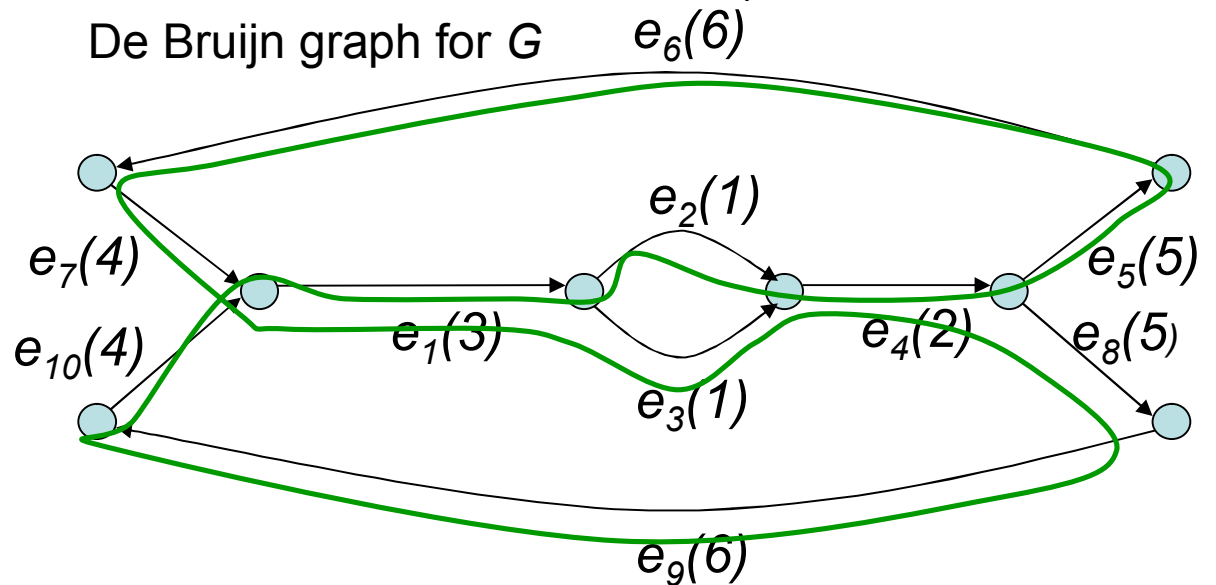
7-bounded paths

$e_1e_2e_4e_5$	$e_1e_3e_4e_8$
$e_2e_4e_5$	$e_3e_4e_8$
$e_4e_5e_6$	$e_4e_8e_9$
$e_5e_6e_7$	$e_8e_9e_{10}$
$e_6e_7e_1$	$e_9e_{10}e_1$
$e_7e_1e_3e_4e_8$	$e_{10}e_1e_2e_4e_5$

Genome G



De Bruijn graph for G

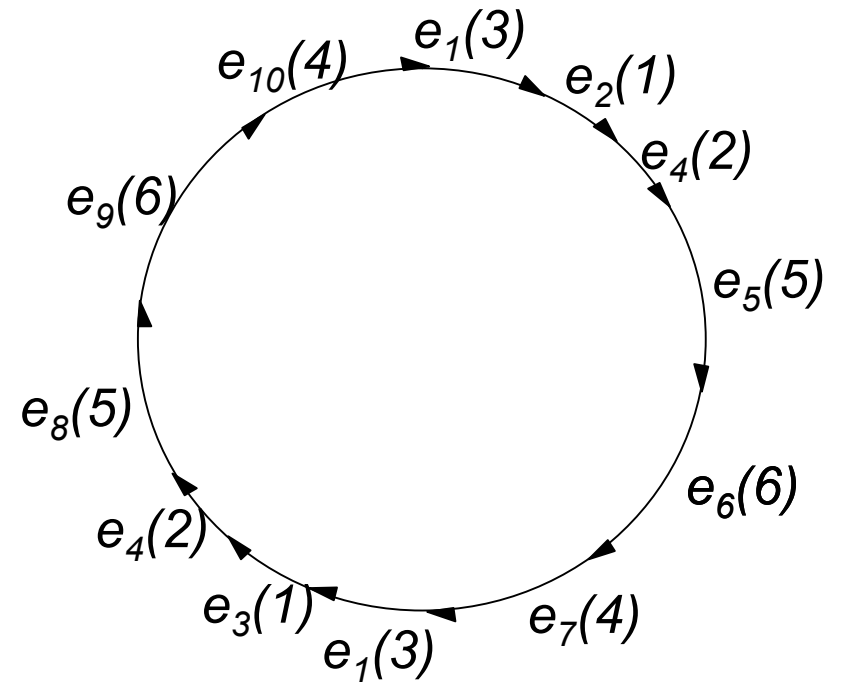


# Pathset assembly

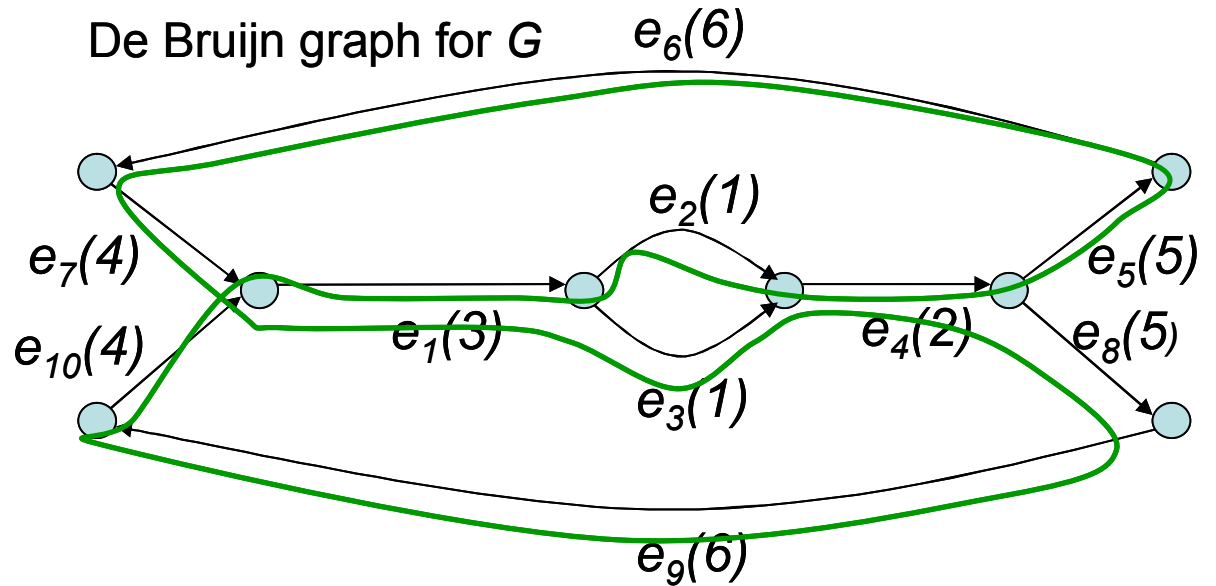
*7-maxpaths*

$e_4e_5e_6$	$e_4e_8e_9$
$e_5e_6e_7$	$e_8e_9e_{10}$
$e_6e_7e_1$	$e_9e_{10}e_1$
$e_7e_1e_3e_4e_8$	$e_{10}e_1e_2e_4e_5$

Genome G



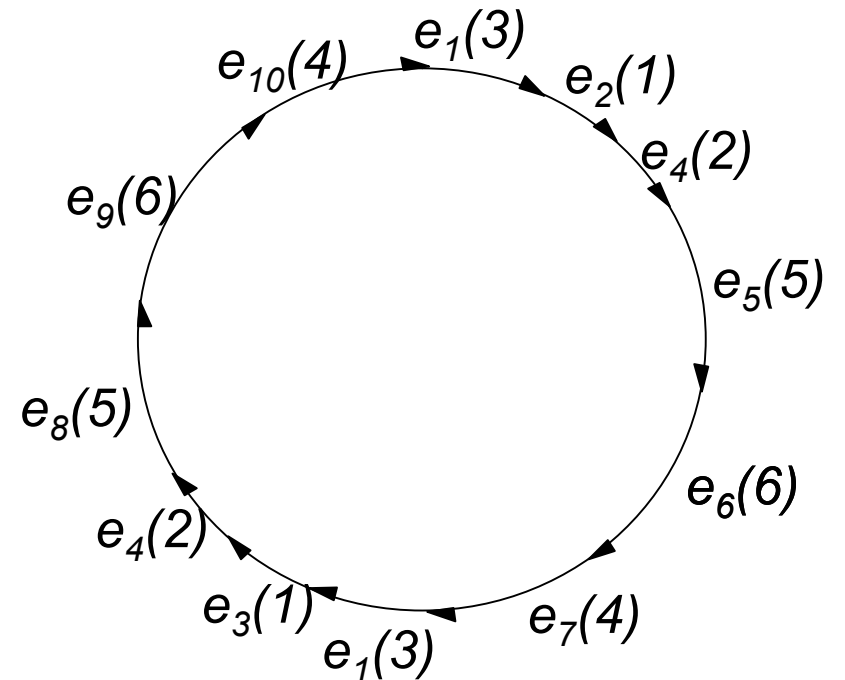
De Bruijn graph for G



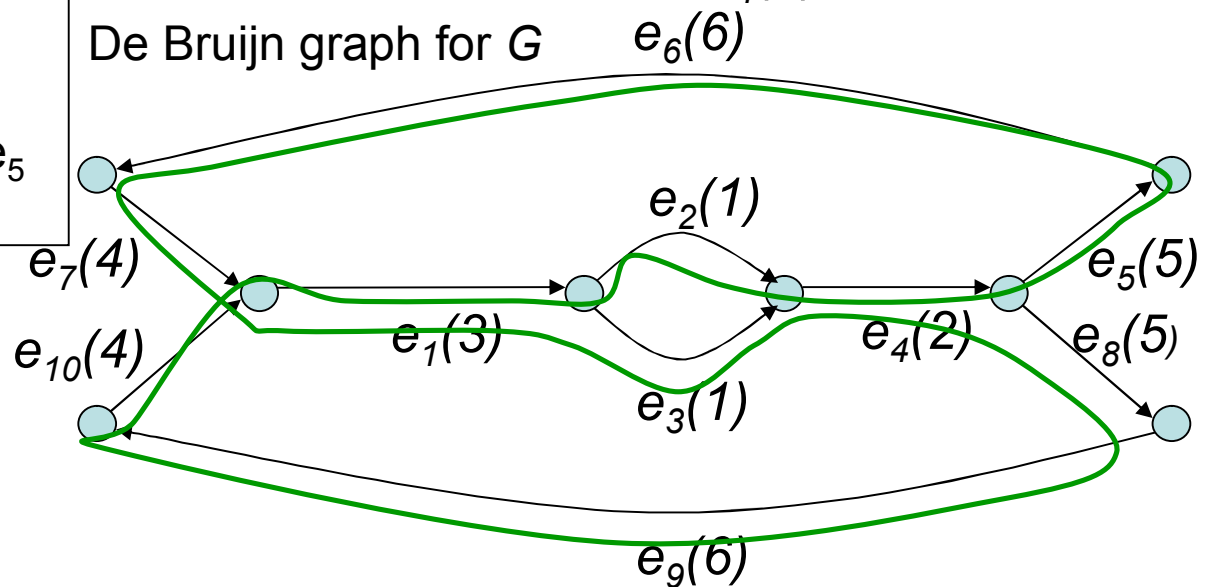
# Pathset assembly

<i>Biedges(a b,d)</i>	<i>7-maxpaths</i>
$(e_4 e_6,7)$	$e_4e_5e_6$
$(e_5 e_7,7)$	$e_5e_6e_7$
$(e_6 e_1,7)$	$e_6e_7e_1$
$(e_7 e_8,7)$	$e_7e_1e_3e_4e_8$
$(e_4 e_9,7)$	$e_4e_8e_9$
$(e_8 e_{10},7)$	$e_8e_9e_{10}$
$(e_9 e_1,7)$	$e_9e_{10}e_1$
$(e_{10} e_5,7)$	$e_{10}e_1e_2e_4e_5$

Genome  $G$



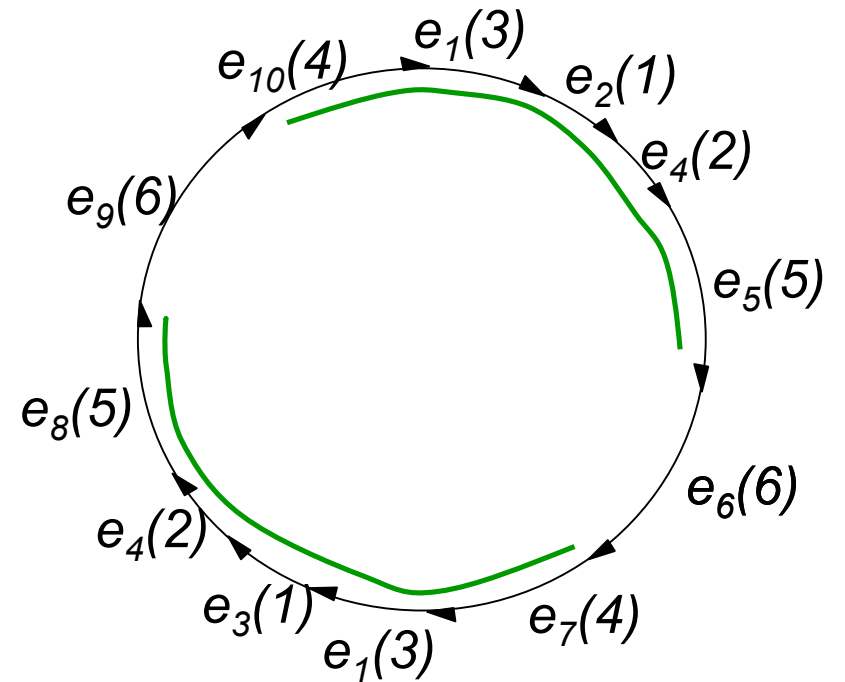
De Bruijn graph for  $G$



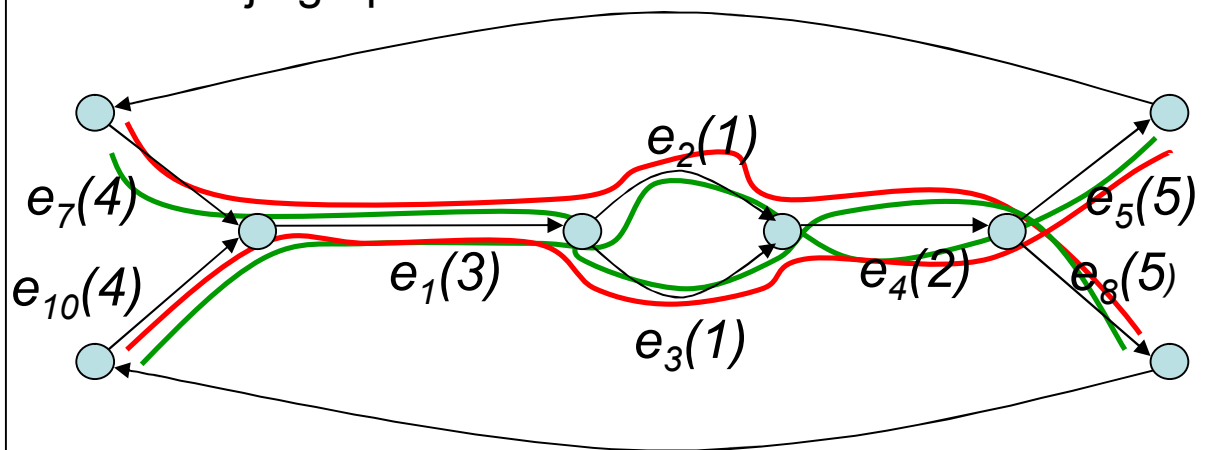
# Pathset assembly

<i>Biedges(a b,d)</i>	<i>Pathsets</i>
$(e_4 e_6,7)$	$\{e_4e_5e_6\}$
$(e_5 e_7,7)$	$\{e_5e_6e_7\}$
$(e_6 e_1,7)$	$\{e_6e_7e_1\}$
$(e_7 e_8,7)$	$\{e_7e_1e_3e_4e_8;$ $e_7e_1e_2e_4e_8\}$
$(e_4 e_9,7)$	$\{e_4e_8e_9\}$
$(e_8 e_{10},7)$	$\{e_8e_9e_{10}\}$
$(e_9 e_1,7)$	$\{e_9e_{10}e_1\}$
$(e_{10} e_5,7)$	$\{e_{10}e_1e_2e_4e_5;$ $e_{10}e_1e_3e_4e_5\}$

Genome G



De Bruijn graph for G





# Pathset assembly

## Pathsets

$\{e_4e_5e_6\}$

$\{e_5e_6e_7\}$

$\{e_6e_7e_1\}$

$\{e_7e_1e_3e_4e_8;$

$e_7e_1e_2e_4e_8\}$

$\{e_4e_8e_9\}$

$\{e_8e_9e_{10}\}$

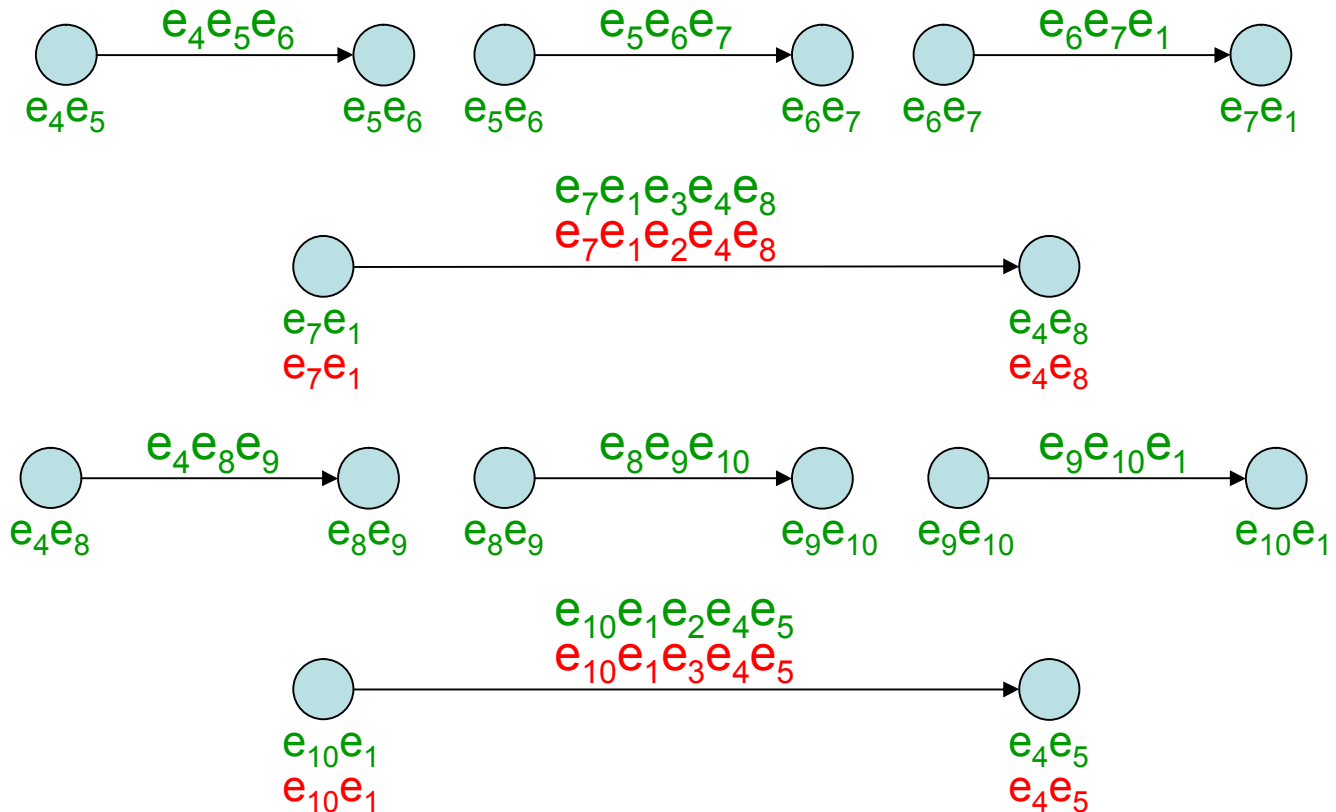
$\{e_9e_{10}e_1\}$

$\{e_{10}e_1e_2e_4e_5;$

$e_{10}e_1e_3e_4e_5\}$

1. For each pathset  $\mathbf{P}$  introduce new vertices  $u, v$  and form edge  $u \rightarrow v$ .

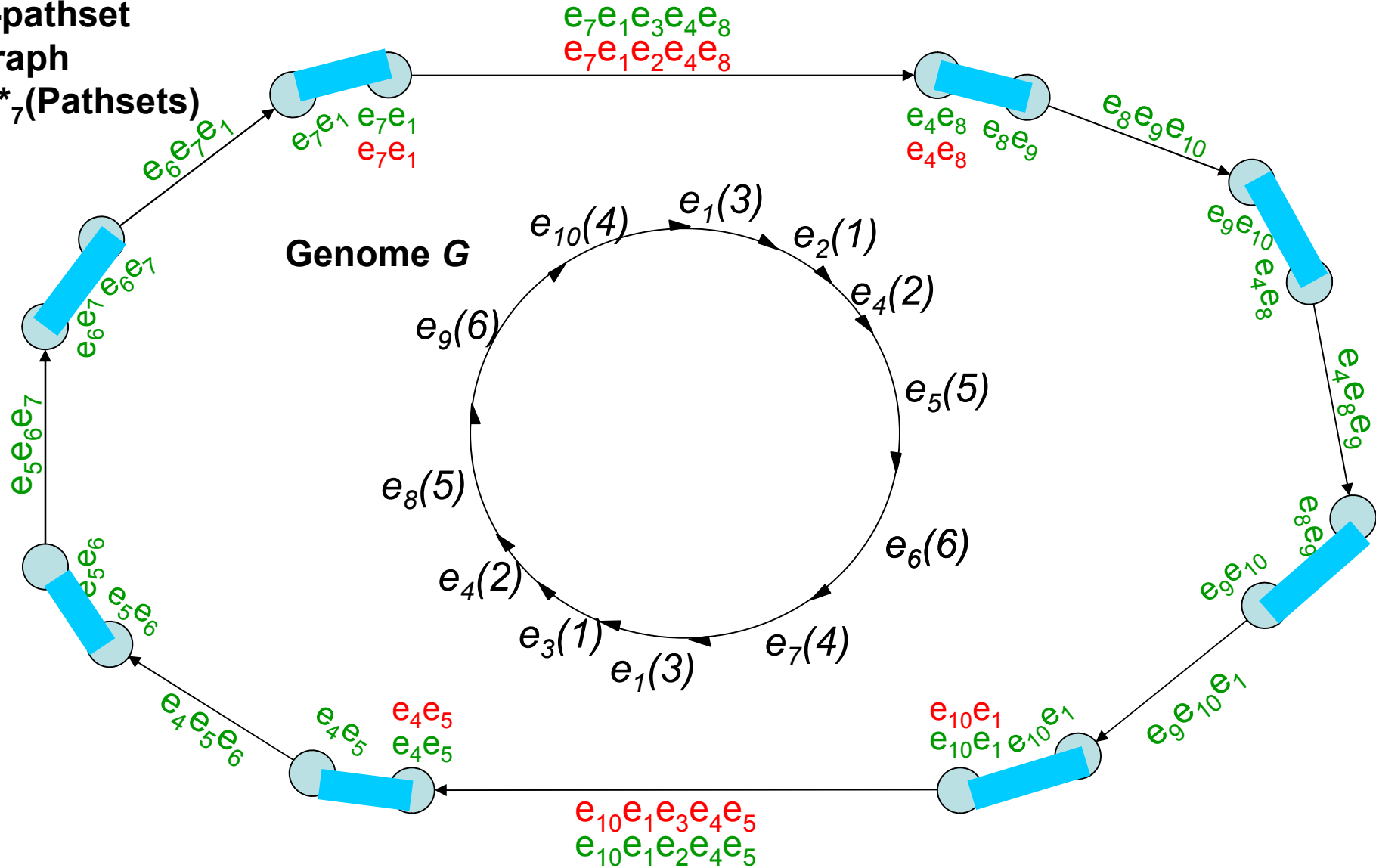
- Label edge by  $\mathbf{P}$ .
- Label  $u$  by set of all minimal long prefixes of paths in  $\mathbf{P}$ .
- Label  $v$  by set of all minimal long suffixes of paths in  $\mathbf{P}$ .



# Pathset assembly

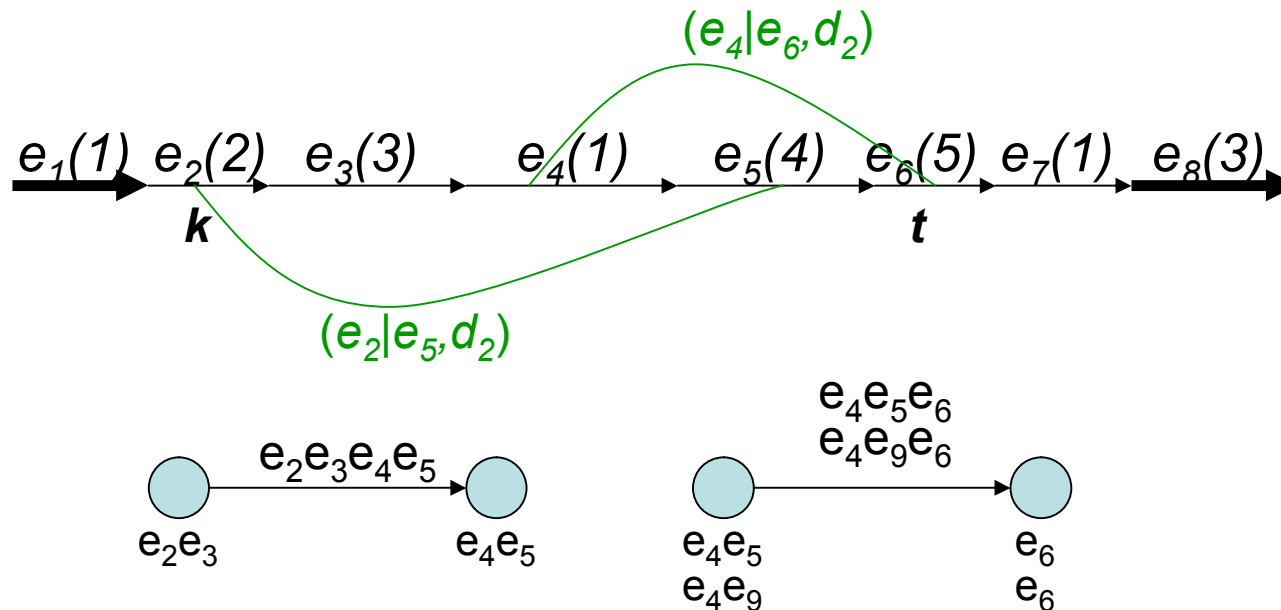
2. Glue vertices together when they share a label.

7-pathset graph  
 $G^*_7(\text{Pathsets})$



# Multiple library pathsets

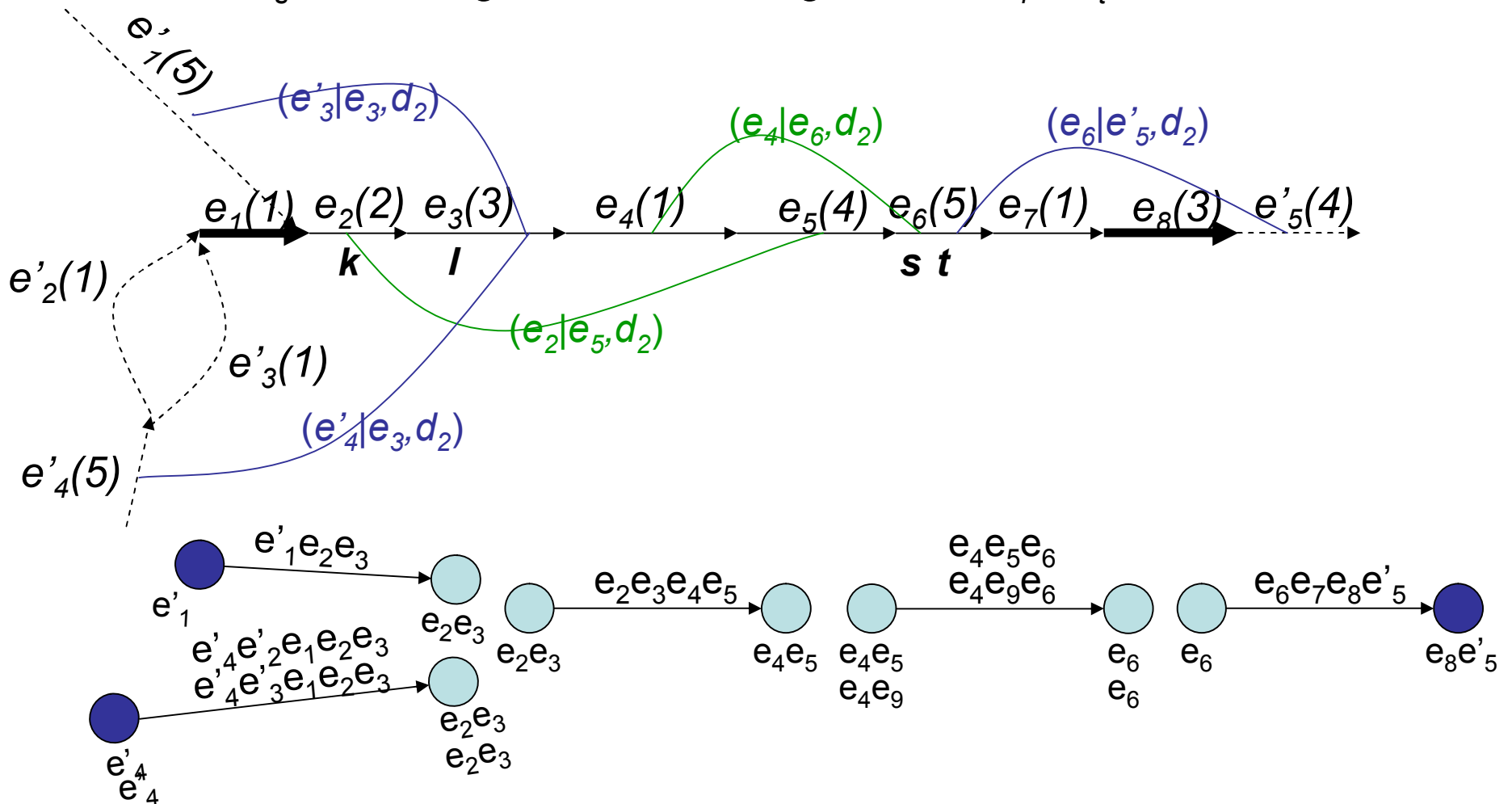
- Given two libraries of biedges  $HB_1$  and  $HB_2$  ( $d_1 > d_2$ )
- Constructed  $Pathset(HB_1)$  for all  $(a|b, d_1) \in HB_1$  and  $Pathset(HB_2)$  for all  $(a|b, d_2) \in HB_2$
- For each path  $\mathbf{p} = e_1, \dots, e_n$  from  $Pathset(a|b, d_1)$  define set of biedges  $HB_{corr}(\mathbf{p})$  and  $Pathset(HB_{corr}(\mathbf{p}))$ :
  - Initiate  $HB_{corr}(\mathbf{p})$  as set of all existing  $(e_i|e_j, d_2) \in HB_2$ ,  $i < j$   
 Find  $\mathbf{k} = \min\{i : (e_i|e_j, d_2) \in HB_{corr}(\mathbf{p})\}$   
 Find  $\mathbf{t} = \max\{j : (e_i|e_j, d_2) \in HB_{corr}(\mathbf{p})\}$



# Multiple library pathsets

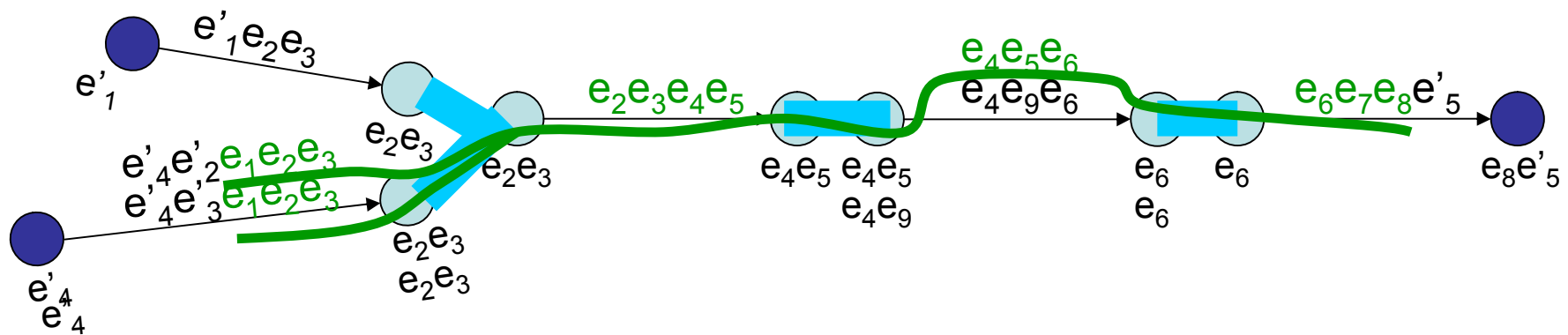
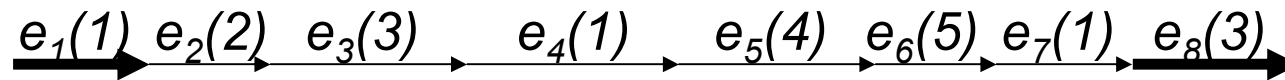
2. If  $e_k \neq e_1$  then add to  $HB_{corr}(\mathbf{p})$  all existing  $(e'|e_l, d_2) \in HB_2$ , where  $e_l$  is last edge of minimal long prefix in  $e_k \dots e_n$

If  $e_t \neq e_n$  then add to  $HB_{corr}(\mathbf{p})$  all existing  $(e_s|e', d_2) \in HB_2$ , where  $e_s$  is first edge of minimal long suffix in  $e_1 \dots e_t$



# Multiple library pathsets

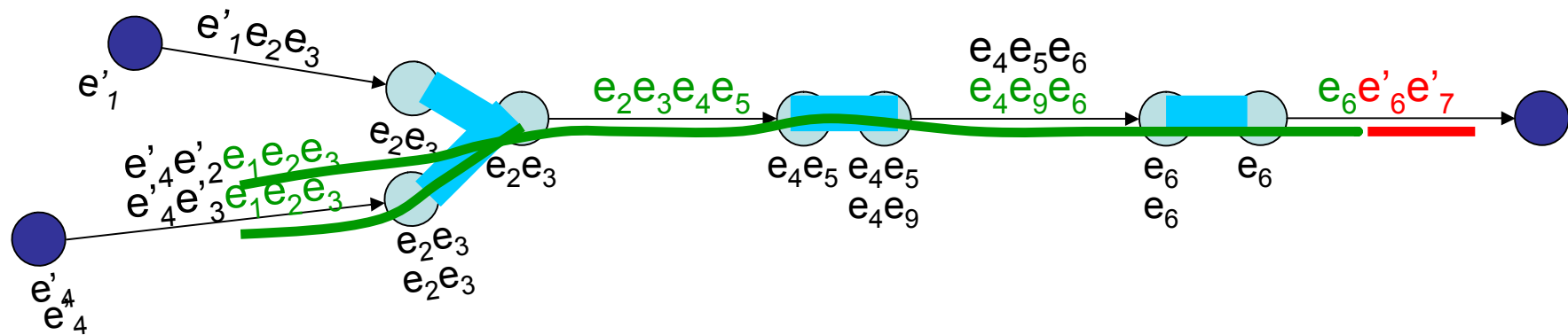
- Construct pathset graph  $G_{d2}^*(HB_{corr}(\mathbf{p}))$   
 If path  $\mathbf{p}$  is genomic, it will be represented in  $G_{d2}^*(HB_{corr})$



# Multiple library pathsets

So, If path  $p$  doesn't exist in  $G_{d_2}^*(HB_{corr}(p))$ , it's not genomic path

$e_1(1) \rightarrow e_2(2) \rightarrow e_3(3) \rightarrow e_4(1) \rightarrow e_9(4) \rightarrow e_6(5) \rightarrow e_7(1) \rightarrow e_8(3)$



- Remove corrupting paths from each  $Pathset(a|b, d_1)$ .  
Removal of such paths reduce some pathsets from  $Pathset(HB_1)$ , it makes pathset graph  $G_{d_1}^*(Pathset(HB_1))$  less ambiguous
- Construct  $G_{d_1}^*(Pathset(HB_1))$  and find Eulerian cycle.