

De Novo Genome Assembly

MPA-PRG: Programming Bioinformatics

Exercise 9

Shortest common superstring

- given a set of strings S , find $SCS(S)$: the shortest string that contains all strings in S as substrings

Example: S : BAA AAB BBA ABA ABB BBB AAA BAB

Concatenation: BAAAABBBAAABAABBBBBBAAABAB

└────────── 24 ─────────┘

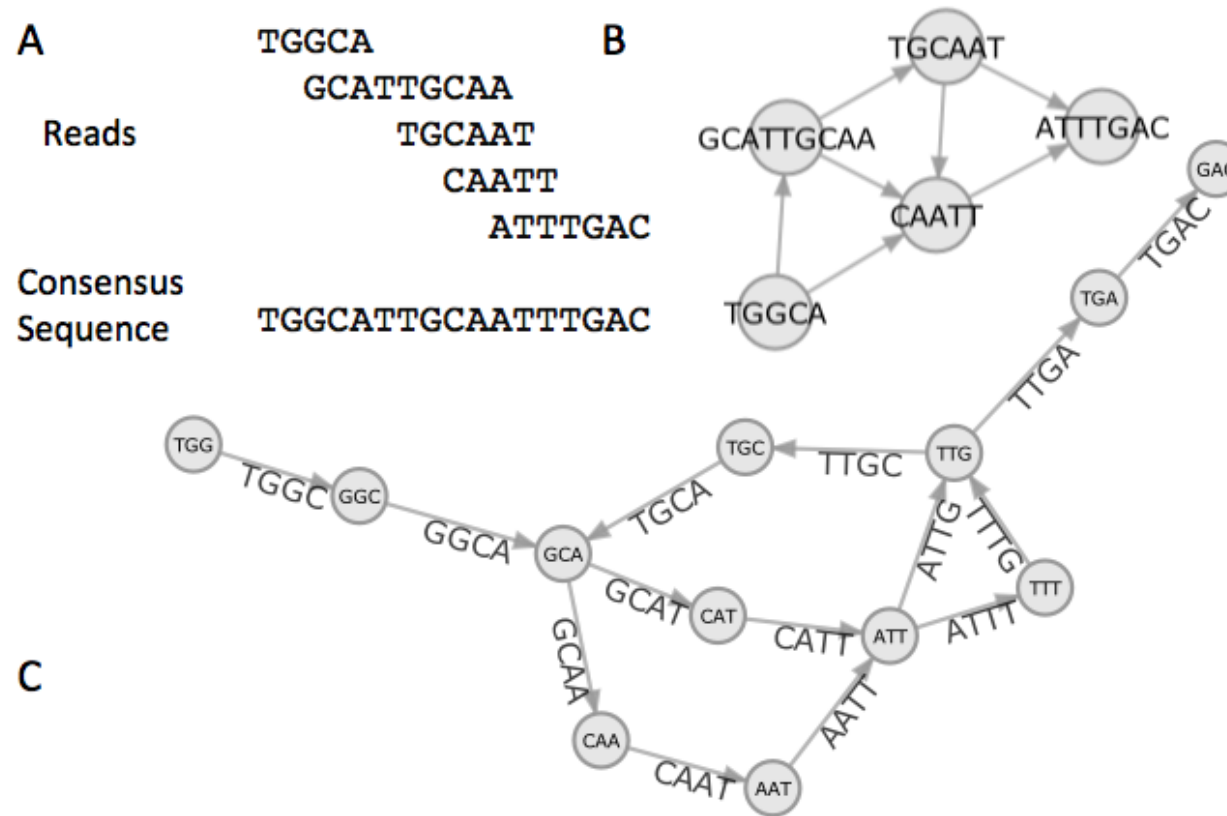
$SCS(S)$: AAABBBABAA

└── 10 ─┘

AAA
AAB
ABB
BBB
BBA
BAB
ABA
BAA

De Novo Assembly

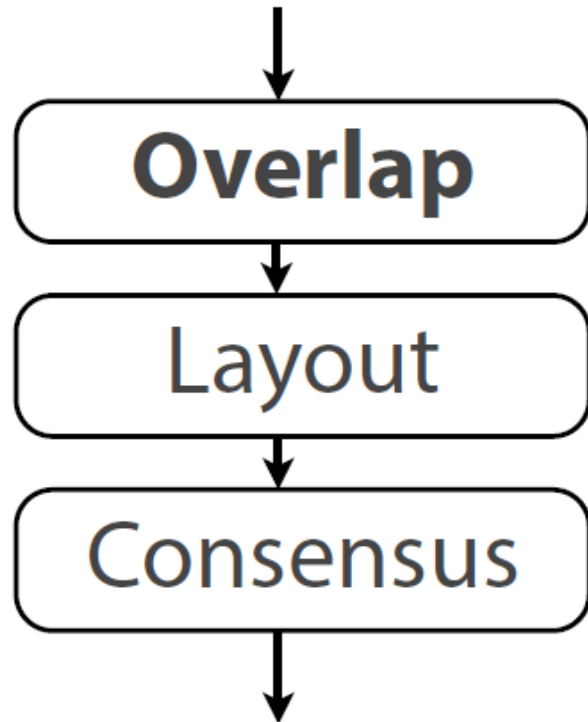
- string methods – greedy algorithms
- graph methods – OLC (overlap layout consensus) × de Bruijn graph



Graph Algorithms

- Overlap Layout Consensus (OLC)
 - reads represent the vertices of the graph and edges the overlaps between them
 - search for a Hamiltonian path (each vertex once)
- de Bruijn graph (DBG)
 - reads represent edges of the graph, and vertices represent overlaps
 - search for a Eulerian walk (each edge once)

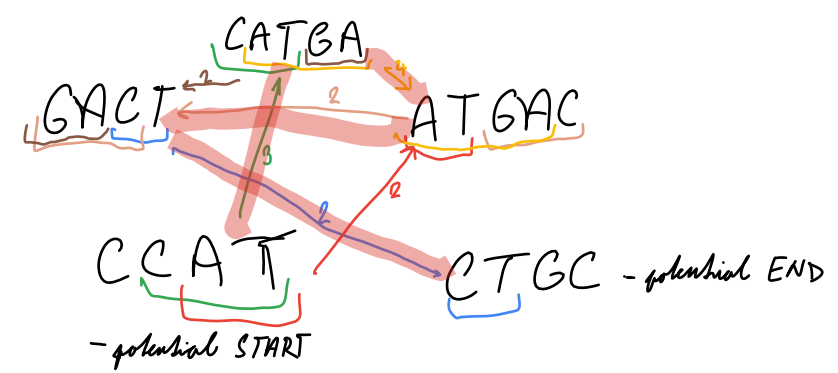
Overlap Layout Consensus



- create an overlap graph
- find an appropriate path through the graph
- select a putative nucleotide sequence

nejdříve se hledá' overlap konce
jednoho se všema oslabníma.

OLC Examples



1. $S = \{GACT, ATGAC, CCAT, CTGC, CATGA\}$

From START \rightarrow biggest path (weight)

2. $S = \{ATG, GACT, CTTA, CATG, ACTTA, TAGCA, GCAT\}$

layout:

```

CCAT
CATGA
ATGAC
GACT
CTGC

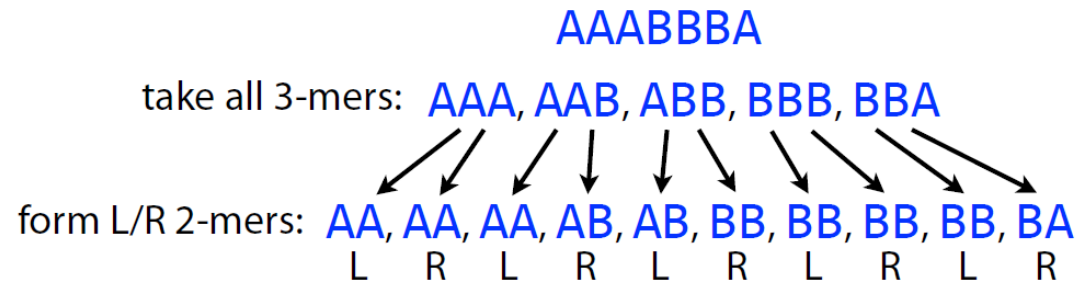
```

Consensus: CCATGACTGC

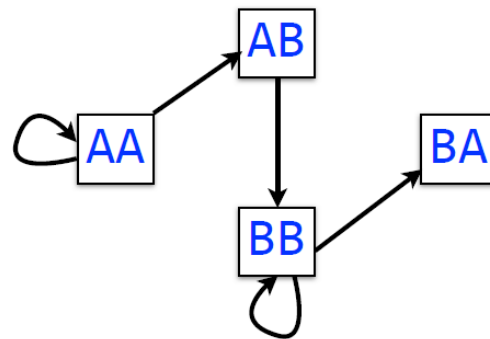
3. $S = \{TAG, AGAG, AGA, GAGT, AGTAGC, GCAG, AGCAG\}$

DBG

Take each length-3 input string and split it into two overlapping substrings of length 2. Call these the *left* and *right* 2-mers.

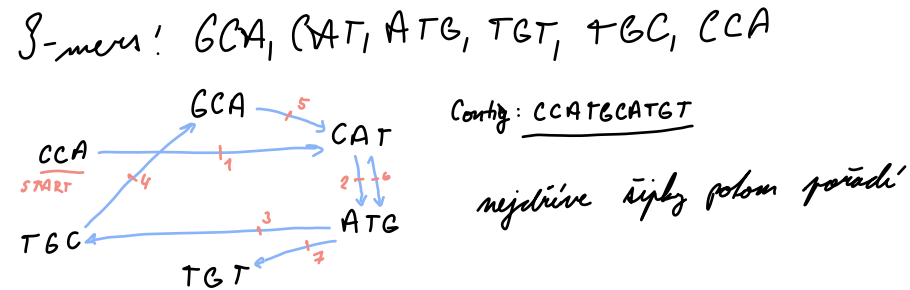


Let 2-mers be nodes in a new graph. Draw a directed edge from each left 2-mer to corresponding right 2-mer:



Each *edge* in this graph corresponds to a length-3 input string

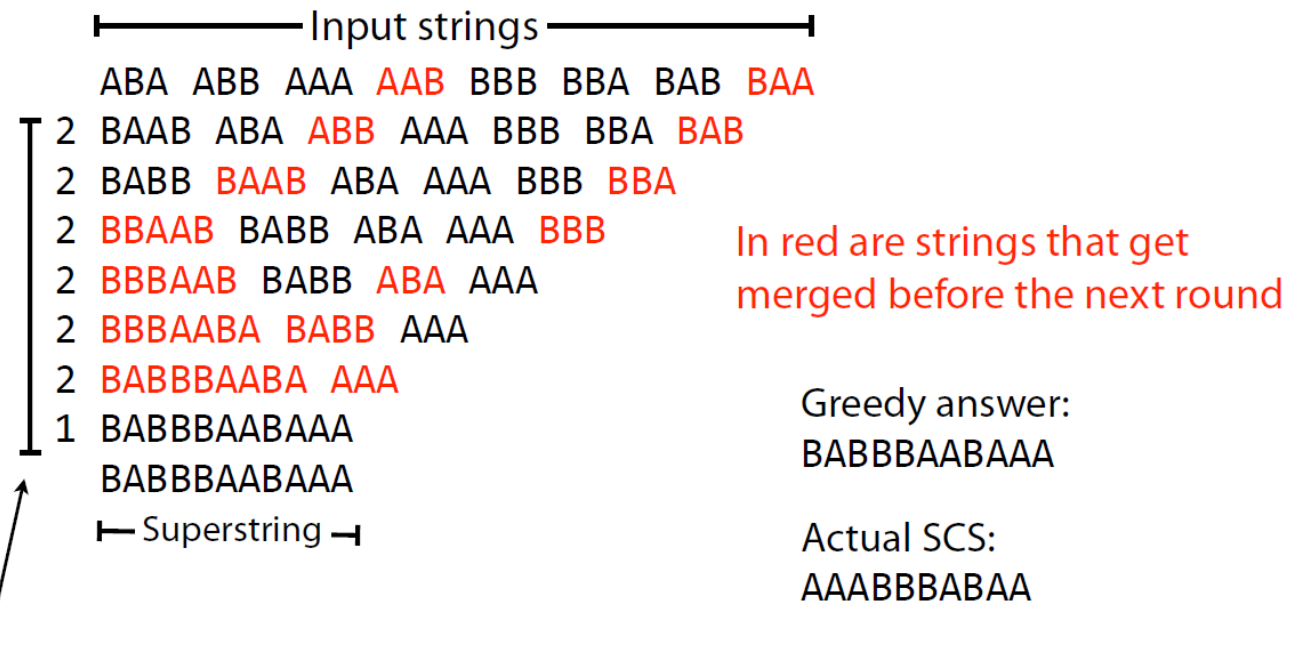
DBG Examples



1. $S = \{GCAT, CATG, ATGT, CATG, ATGC, TGCA, CCAT\}$
2. $S = \{ACGT, AATC, ATCA, AGTA, ACGC, CGTA, CAGT, GTAA, GTAC, TAAT, TCAG, TACG\}$
3. $S = \{ACGTA, TAAT, ATCAGT, GTA, TACG, CGT\}$

String Methods

- greedy algorithms
 - in each step merge the two reads with the largest overlap



Rounds of merging, one merge per line.

Number in first column = length of overlap merged before that round.

Greedy SCS Examples

1. $S = \{\text{CATGC}, \text{CTAAGT}, \text{GCTA}, \text{TTCA}, \text{ATGCATC}\}$

2. $S = \{\text{ATC}, \text{TCAGAG}, \text{ATG}, \text{AGCCAT}, \text{TGCAT}\}$

	<u>CATGC</u>	<u>CTAAGT</u>	<u>GCTA</u>	<u>TTCA</u>	<u>ATGCATC</u>
<u>CATGC</u>	—	1 ^C	2 ^{GC}	0	4 ^{ATGC}
<u>CTAAGT</u>	0	—	0	1 ^T	0
<u>GCTA</u>	2 ^{GC}	3 ^{CTA}	—	0	1 ^A
<u>TTCA</u>	2 ^{CA}	0	0	—	1 ^A
<u>ATGCATC</u>	1 ^C	1 ^C	0	0	—

Task

- In R, implement a function `GreedySuperstring()` according to the pseudocode.
- Input:
 - `S` a `DNAStringSet` of strings (reads)
- Output:
 - `S` a `DNAStringSet` of the shortest common superstring (contig)

```
GreedySuperstring(S)
1  while length of S > 1
2    overlapMat ← OverlapMatrix(S)
3    if max(overlapMat) = 0
4      return S
5    else
6      seq1, seq2 ← Two sequences from S with the longest overlap
7      Merge seq1 and seq2 and add the new sequence to S
8      Remove seq1 and seq2 from S
9  return S
```

Hint: Create also functions:

- `Overlap()` to calculate overlap between two sequences
- `OverlapMatrix()` to create a matrix of overlaps among all sequences in `S`