

CS267 Assignment 3:

Parallelize Graph Algorithms for de Novo Genome Assembly

Spring 2016

2

Problem statement

- **Input:** A set of unique k-mers and their corresponding extensions.
 - k-mers are sequences of length k (alphabet is A/C/G/T).
 - An extension is a simple symbol (A/C/G/T/F).
 - The input k-mers form a de Bruijn graph, a special graph that is used to represent overlaps between sequences of symbols.
- **Output:** A set of contigs, i.e. connected components in the input de Bruijn graph.

3

Example

- **Input:** A set of unique k-mers and their corresponding extensions.
- Example for k = 3:
- Format:

k-mer	forward extension , backward extension
AAC	CF
ATC	TG
ACC	GA
TGA	FC
GAT	CF
AAT	GF
ATG	CA
TCT	GA
CCG	FA
CTG	AT
TGC	FA

4

Example

- **Input:** A set of unique k-mers and their corresponding extensions.
- The input corresponds to a de Bruijn graph.
- Example for k = 3:

AAC	CF
ATC	TG
ACC	GA
TGA	FC
GAT	CF
AAT	GF
ATG	CA
TCT	GA
CCG	FA
CTG	AT
TGC	FA

Sequence of k-mers :

Sequence of k-mers :

Sequence of k-mers :

5

Example

- **Input:** A set of unique k-mers and their corresponding extensions.
 - The input corresponds to a de Bruijn graph.
- Example for k = 3:

AAC ATC ACC TGA GAT AAT ATG TCT CCG CTG TGC	CF TG GA FC CF GF CA GA FA AT FA	Sequence of k-mers : 	Sequence of k-mers :
		Sequence of k-mers : 	Sequence of k-mers :

k-mers with "F" as an extension are start vertices

6

Example

- **Input:** A set of unique k-mers and their corresponding extensions.
 - The input corresponds to a de Bruijn graph.
- Example for k = 3:

AAC ATC ACC TGA GAT AAT ATG TCT CCG CTG TGC	CF TG GA FC CF GF CA GA FA AT FA	Sequence of k-mers : 	Sequence of k-mers :
		Sequence of k-mers : 	Sequence of k-mers :

- Consider k-mer: TCT
- Concatenate last k-1 bases (CT) and forward extension (G) => CTG (following vertex)
- Concatenate backward extension (A) and first k-1 bases (TC) =>ATC (preceding vertex)
- The graph is undirected, we can visit a vertex from both directions.

7

Example

- **Input:** A set of unique k-mers and their corresponding extensions.
 - The input corresponds to a de Bruijn graph.
- Example for k = 3:

AAC ATC ACC TGA GAT AAT ATG TCT CCG CTG TGC	CF TG GA FC CF GF CA GA FA AT FA	Sequence of k-mers : 	Contig : AACCG
		Contig : GATCTGA 	Contig : AATGC

- **Output:** A set of contigs or equivalently the connected components in the de Bruijn graph

8

Compact graph representation: hash table

- The vertices are keys
- The edges (neighboring vertices) are represented with a two-letter value

AAC ATC ACC TGA GAT AAT ATG TCT CCG CTG TGC	CF TG GA FC CF GF CA GA FA AT FA	buckets 	entries
---	--	-------------	-------------

Serial algorithm

9

Algorithm 1 De Bruijn Graph Construction And Traversal

```

1: Input: A set of  $k$ -mers and their corresponding forward and backward extensions
2: Output: A set of contigs
3: /* Initialization */
4:  $hashTable \leftarrow CREATEHASHTABLE()$ 
5:  $startNodesList \leftarrow CREATEEMPTYLIST()$ 
6:
7: /* De Bruijn Graph Construction */
8: for each ( $k$ -mer, forwardExt, backwardExt) in input do
9:    $ADDKMERHASHTABLE(hashTable, (k$ -mer, forwardExt, backwardExt))
10:  if (backwardExt is F) then
11:     $ADDKMERLIST(startNodesList, (k$ -mer, forwardExt))
12:  end if
13: end for
14:
15: /* De Bruijn Graph Traversal */
16: for each ( $k$ -mer, forwardExt) in startNodesList do
17:    $currentContig \leftarrow CREATENEWSEQUENCE(k$ -mer)
18:    $currentForwardExtension \leftarrow forwardExt$ 
19:   while ( $currentForwardExtension$  is not F) do
20:      $ADDBASETOSEQUENCE(currentForwardExtension, currentContig)$ 
21:      $currentKmer \leftarrow LASTKBASES(currentContig)$ 
22:      $currentForwardExtension \leftarrow LOOKUP(hashTable, currentKmer)$ 
23:   end while
24:    $STORECONTIG(currentContig)$ 
25: end for
    
```

Graph construction

10

- The vertices are keys
- The edges (neighboring vertices) are represented with a two-letter value

The diagram illustrates the mapping from buckets to entries. Buckets are keys (3-letter strings) and entries are pairs of (key, forward/backward extensions). The buckets are: ATC, AAC, TGA, GAT, AAT, TCT, CCG, CTG. The entries are: (ATC, T/G), (AAC, C/F), (TGA, F/C), (GAT, C/F), (AAT, G/F), (TCT, G/A), (CCG, F/A), (CTG, A/T). Edges connect buckets to their corresponding entries.

Graph traversal

11

- We pick a start vertex and we initiate a contig.

Contig: A A T

The graph shows nodes (keys) and edges (two-letter values). The path AAT is highlighted with a red dashed box. The nodes are: AAC, ATC, ACC, TGA, GAT, AAT, ATG, TCT, CCG, CTG, TGA, AAT, ATG, TGC.

Graph traversal

12

- We add the forward extension to the contig.

Contig: A A T G

The graph shows nodes (keys) and edges (two-letter values). The path AATG is highlighted with a red dashed box. The nodes are: AAC, ATC, ACC, TGA, GAT, AAT, ATG, TCT, CCG, CTG, TGA, AAT, ATG, TGC.

13

Graph traversal

- We take the last k bases of the contig and look them up in the hash table.

Contig: **AATG**

AAC	CF
ATC	TG
ACC	GA
TGA	FC
GAT	CF
AAT	GF
ATG	CA
TCT	GA
CCG	FA
CTG	AT
TGC	FA

14

Graph traversal

- We add the new forward extension to the contig.

Contig: **AATGC**

AAC	CF
ATC	TG
ACC	GA
TGA	FC
GAT	CF
AAT	GF
ATG	CA
TCT	GA
CCG	FA
CTG	AT
TGC	FA

15

Graph traversal

- We take the last k bases of the contig and look them up in the hash table.

Contig: **AATGC**

AAC	CF
ATC	TG
ACC	GA
TGA	FC
GAT	CF
AAT	GF
ATG	CA
TCT	GA
CCG	FA
CTG	AT
TGC	FA

16

Graph traversal

- We terminate the current contig since the forward extension is an "F".

Contig: **AATGC**

AAC	CF
ATC	TG
ACC	GA
TGA	FC
GAT	CF
AAT	GF
ATG	CA
TCT	GA
CCG	FA
CTG	AT
TGC	FA

