

## CS267 Assignment 3:

# Parallelize Graph Algorithms for de Novo Genome Assembly

Spring 2015

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

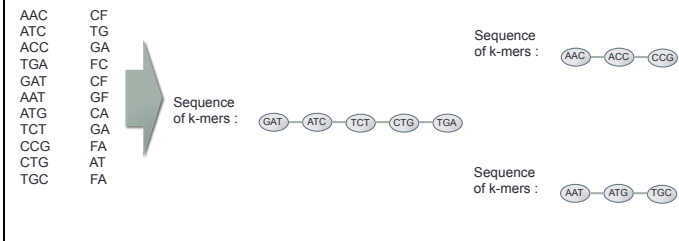
- **Input:** A set of unique k-mers and their corresponding extensions.
- Example for k = 3:
- Format:

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:



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 <table border="1" style="width: 100%; border-collapse: collapse;"> <tr> <td>key: ATC</td> <td>forw_ext: T back_ext: G</td> <td>key: ACC</td> <td>forw_ext: G back_ext: A</td> </tr> <tr> <td>key: AAC</td> <td>forw_ext: C back_ext: F</td> <td></td> <td></td> </tr> <tr> <td>key: TGA</td> <td>forw_ext: F back_ext: C</td> <td></td> <td></td> </tr> <tr> <td>key: GAT</td> <td>forw_ext: C back_ext: F</td> <td>key: ATG</td> <td>forw_ext: C back_ext: A</td> </tr> <tr> <td>key: AAT</td> <td>forw_ext: G back_ext: F</td> <td></td> <td></td> </tr> <tr> <td>key: TCT</td> <td>forw_ext: G back_ext: A</td> <td></td> <td></td> </tr> <tr> <td>key: CCG</td> <td>forw_ext: F back_ext: A</td> <td></td> <td></td> </tr> <tr> <td>key: CTG</td> <td>forw_ext: A back_ext: T</td> <td>key: TGC</td> <td>forw_ext: F back_ext: A</td> </tr> </table>	key: ATC	forw_ext: T back_ext: G	key: ACC	forw_ext: G back_ext: A	key: AAC	forw_ext: C back_ext: F			key: TGA	forw_ext: F back_ext: C			key: GAT	forw_ext: C back_ext: F	key: ATG	forw_ext: C back_ext: A	key: AAT	forw_ext: G back_ext: F			key: TCT	forw_ext: G back_ext: A			key: CCG	forw_ext: F back_ext: A			key: CTG	forw_ext: A back_ext: T	key: TGC	forw_ext: F back_ext: A
key: ATC	forw_ext: T back_ext: G	key: ACC	forw_ext: G back_ext: A																																
key: AAC	forw_ext: C back_ext: F																																		
key: TGA	forw_ext: F back_ext: C																																		
key: GAT	forw_ext: C back_ext: F	key: ATG	forw_ext: C back_ext: A																																
key: AAT	forw_ext: G back_ext: F																																		
key: TCT	forw_ext: G back_ext: A																																		
key: CCG	forw_ext: F back_ext: A																																		
key: CTG	forw_ext: A back_ext: T	key: TGC	forw_ext: F back_ext: A																																

9

## Serial algorithm

---

**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 ← CREATEHASHTABLE()
5: startNodesList ← CREATEEMPTYLIST()
6:
7: /* De Bruijn Graph Construction */
8: for each ( $k$ -mer, forwardExt, backwardExt) in input do
9:   ADDKMERTOHASHTABLE(hashTable, ( $k$ -mer, forwardExt, backwardExt))
10:  if (backwardExt is F) then
11:    ADDKMERTOLIST(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 ← CREATENEWSEQUENCE( $k$ -mer)
18:   currentForwardExtension ← forwardExt
19:   while (currentForwardExtension is not F) do
20:     ADDBASETOSEQUENCE(currentForwardExtension, currentContig)
21:     currentKmer ← LASTKBASES(currentContig)
22:     currentForwardExtension ← LOOKUP(hashTable, currentKmer)
23:   end while
24:   STORECONTIG(currentContig)
25: end for
    
```

10

## Graph construction

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

key	forw_ext	back_ext
ATC	T	G
AAC	C	F
TGA	F	C
GAT	C	F
AAT	G	F
TCT	G	A
CCG	F	A
CTG	A	T

11

## Graph traversal

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

**Contig: A A T**

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

12

## Graph traversal

- We add the forward extension to the contig.

**Contig: A A T G**

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

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
<b>ATG</b>	<b>CA</b>
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
<b>ATG</b>	<b>CA</b>
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
<b>TGC</b>	<b>FA</b>

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
<b>TGC</b>	<b>FA</b>

