# EVALUATION AND BENCHMARKING OF A NEW SCAFFOLDING METHODOLOGY

## ALEXANDRINA BODRUG

SUPERVISORS: PR. RUMEN ANDONOV & DR. DOMINIQUE LAVENIER

### UNIVERSITY RENNES 1
BIOINFORMATICS AND GENOMICS MASTER

JUNE 25, 2015

# Overview

SCAFFOLDING

ALEXANDRINA
BODRUG

Context
Some definitions
Order and orient
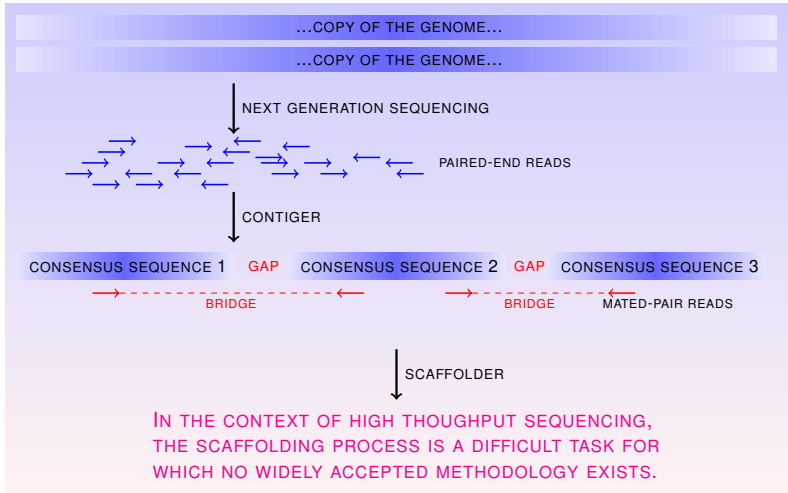Challenging problem

GST features
Scripting for the GST
GST modeled graph
Expected solution
Scaffolding solutions

Benchmarking
workflow for
the GST

Results
Large repeats
Short repeats

Perspectives

# CONTEXT

...COPY OF THE GENOME...

...COPY OF THE GENOME...

NEXT GENERATION SEQUENCING

PAIRED-END READS

CONTIGER

CONSENSUS SEQUENCE 1  GAP  CONSENSUS SEQUENCE 2  GAP  CONSENSUS SEQUENCE 3

BRIDGE        BRIDGE        MATED-PAIR READS

SCAFFOLDER

IN THE CONTEXT OF HIGH THOUGHPUT SEQUENCING,
THE SCAFFOLDING PROCESS IS A DIFFICULT TASK FOR
WHICH NO WIDELY ACCEPTED METHODOLOGY EXISTS.

# SOME DEFINITIONS

"The *Contig Scaffolding Problem* is to order and orientate the given contigs in a manner that is consistent with as many mate-pairs as possible".
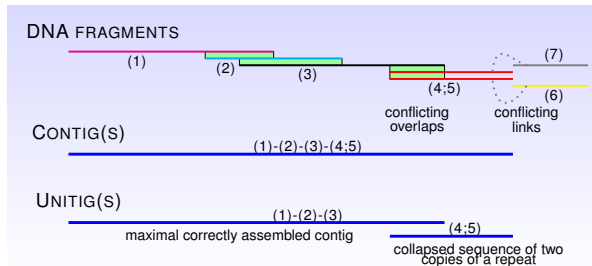
Hudson *et al.* 2002

Genome is fragmented, extremities are sequenced ($\mapsto$ reads) . . .



. . . reads are assembled into consensus sequences.
UNITIGS ARE HIGH CONFIDENCE CONTIGS.

Mated-pair read $\mapsto$ bridge between contigs
Several correctly mapped reads $\mapsto$ link between contigs
High-confidence overlap $\mapsto$ link between contigs
All linkage information $\mapsto$ order and orient contigs

CONFLICTING LINKS CAN AND WILL EXIST.

# CHALLENGING PROBLEM

What would you do in these situations?

## SEVERAL TOOLS MODEL THE PROBLEM DIFFERENTLY.

$\rightarrow$ common features:

- modeling of the scaffolding problem as a graph
- use of unitigs instead of contigs to better compute coverage
- use of unitig coverages to duplicate nodes representing unitigs
- unitig orientations represented by separate nodes

$\rightarrow$ differences:

- **weighted path model** focuses solely on order and orientation
- **distance based model** incorporates link length information
- **flow model** accepts intervals for unitig coverage and link length

# HANDLING THE MODELED GRAPHS

Automated ways to control the input data and validate the scaffolding solution:

- a script to visualize input data and GST solutions: `graph_generator.py`
- a script to inspect the features of the modeled input graph: `graph_inspector.py`
- a script to automatically detect correctly solved instances: `graph_comparator.py`

# GST MODELED GRAPH

SCAFFOLDING

ALEXANDRINA
BODRUG

Context
Some definitions
Order and orient
Challenging problem

GST features
Scripting for the GST
**GST modeled graph**
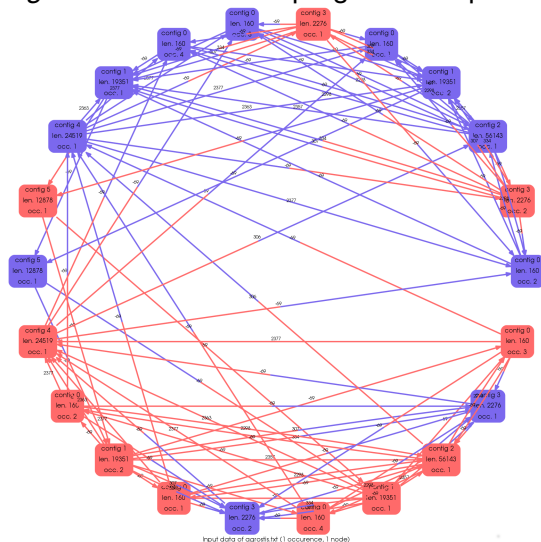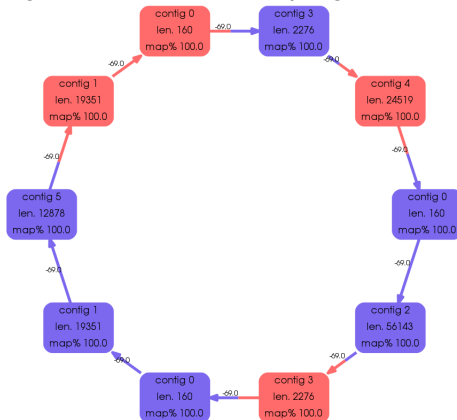Expected solution
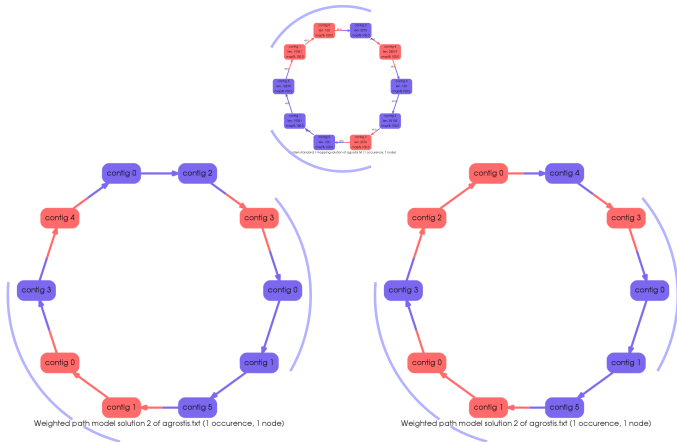Scaffolding solutions

Benchmarking
workflow for
the GST

Results
Large repeats
Short repeats

Perspectives

## *Agrostis stolonifera* chpl. genome input data graph



Input data of agrostis.txt (1 occurence, 1 node)

$$\sum nodes = \sum cov. \times 2$$

| unitig | len | cov |
|--------|-------|--------|
| 1 | 19351 | [2, 2] |
| 0 | 160 | [2, 4] |
| 3 | 2276 | [2, 2] |
| 2 | 56143 | [1, 1] |
| 5 | 12878 | [1, 1] |
| 4 | 24519 | [1, 1] |

*Agrostis stolonifera* chpl. genome expected solution



Golden standard / mapping solution of agrostis.txt (1 occurence, 1 node)

→USE UNITIGS THE RIGHT NUM-
BER OF TIMES IN THE CORRECT
ORIENTATION
→ORDER UNITIGS TO OBTAIN AN
UNINTERRUPTED CIRCULAR PATH

| unitig | orient. | occ. |
|--------|---------|------|
| 1 | reverse | 1 |
| 1 | forward | 1 |
| 3 | reverse | 1 |
| 3 | forward | 1 |
| 0 | reverse | 1 |
| 0 | forward | 2 |

# Scaffolding solutions

## *Agrostis stolonifera* chpl. genome scaffolding solutions

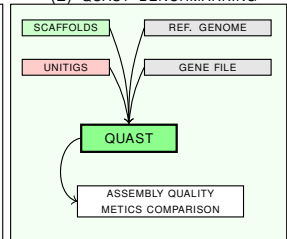### UNITIGS FORMING THE INVERTED REPEATED SEQUENCE OF THE CHLOROPLASTIC GENOME ARE DUPLICATED.



Weighted path model solution 2 of agrostis.txt (1 occurence, 1 node)

# BENCHMARKING WORKFLOW FOR THE GST

SOLUTIONS FOUND WITH THE GENSCALE TOOLS ARE BENCHMARKED
AGAINST THE SSPACE PUBLISHED SCAFFOLDER.

# RESULTS
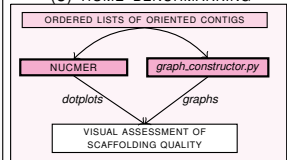
Using the benchmark workflow the following conclusions
were drawn:

- Genomes with big repeated regions are solved a lot
  better with GSTs than with SSPACE
- Small repeats are very challenging to scaffold because
  too many conflicting links exist and GST can not take a
  decision or is too slow
- The GST models processing the link sequence length
  information perform worse than those focusing only on
  ordering and orientating

## EXCELLENT RESULTS ARE OBTAINED FOR DATA SETS WITH LARGE REPEATS AND A SMALL NUMBER OF UNITIGS.

# SHORT REPEATS - BACTERIAL GENOMES

SMALL REPEATS ARE PROBLEMATIC STARTING FROM THE UNITIG BUILDING STEP.
THE WOLBACHIA ENDOSYMBIONT ORGANISM POSSESSES 444 UNITIGS AND ONLY
138 ARE LONGER THAN 1000 BASE PAIRS.

### DEVELOP - TEST - BENCHMARK

- Find strategies which solve more challenging data
  $\rightarrow$ flow model in development
- Benchmark against other tools trying to solve repeated regions
- Test the GST with real data
- Test the GST with other genome sequencing data types

SCAFFOLDING: SAFETY COMES FIRST