

## 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

# 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 24, 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

## 1 Context

- Some definitions
- Order and orient
- Challenging problem

## 2 GST features

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

## 3 Benchmarking workflow for the GST

## 4 Results

- Large repeats
- Short repeats

## 5 Perspectives

# CONTEXT

## SCAFFOLDING

ALEXANDRINA  
BODRUG

### Context

Some definitions  
Order and orient  
Challenging problem

### GST features

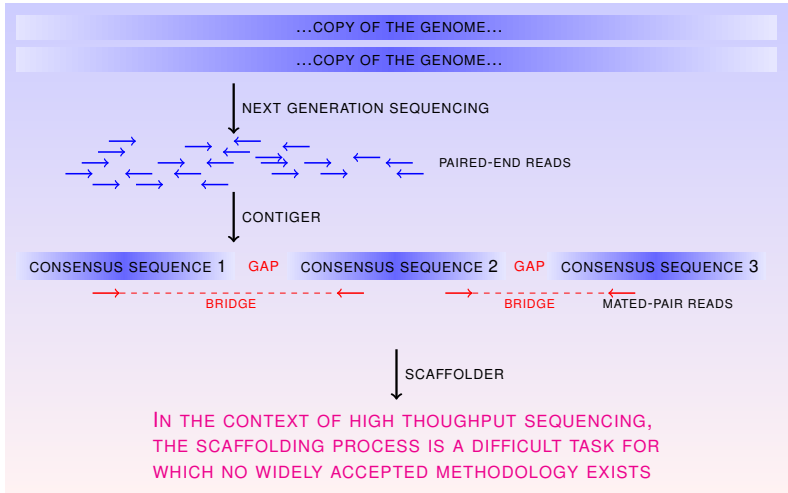
Scripting for the GST  
GST modeled graph  
Expected solution  
Scaffolding solutions

### Benchmarking workflow for the GST

### Results

Large repeats  
Short repeats

### Perspectives



# SOME DEFINITIONS

## 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

"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

# SOME DEFINITIONS

## 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

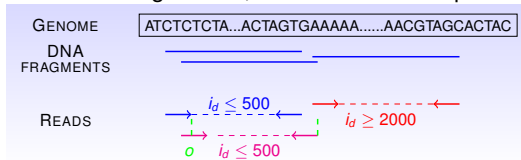
### Results

Large repeats

Short repeats

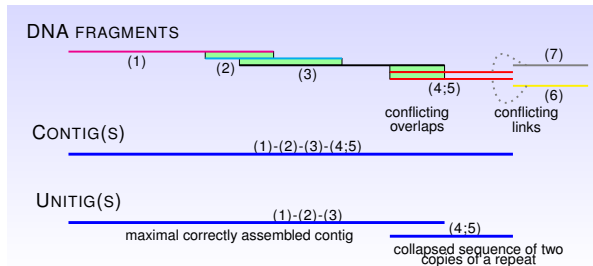
### Perspectives

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



... reads are assembled into consensus sequences.

**UNITIGS ARE HIGH CONFIDENCE CONTIGS.**



# ORDER AND ORIENT

## 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

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

## 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

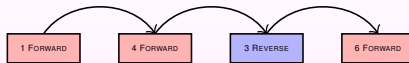
## Results

Large repeats

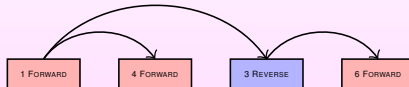
Short repeats

## Perspectives

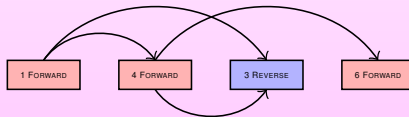
What would you do in these situations?



EASY PATH



LINK COVERAGE?



HEURISTICS...

# GENSCALE SCAFFOLDING TOOLS FEATURES

## 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

## SEVERAL TOOLS MODEL THE PROBLEM DIFFERENTLY.

→ common features:

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

→ 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

## 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

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`





# SCAFFOLDING SOLUTIONS

## 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

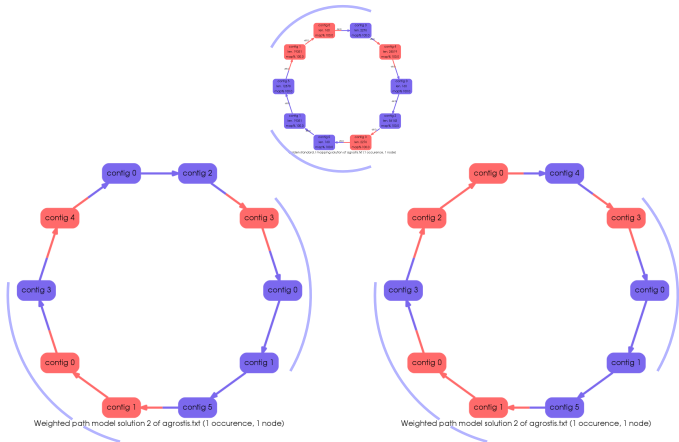
### Results

Large repeats  
Short repeats

### Perspectives

## *Agrostis stolonifera* chl. genome scaffolding solutions

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



# BENCHMARKING WORKFLOW FOR THE GST

## SCAFFOLDING

ALEXANDRINA  
BODRUG

### Context

Some definitions  
Order and orient  
Challenging problem

### GST features

Scripting for the GST  
GST modeled graph  
Expected solution  
Scaffolding solutions

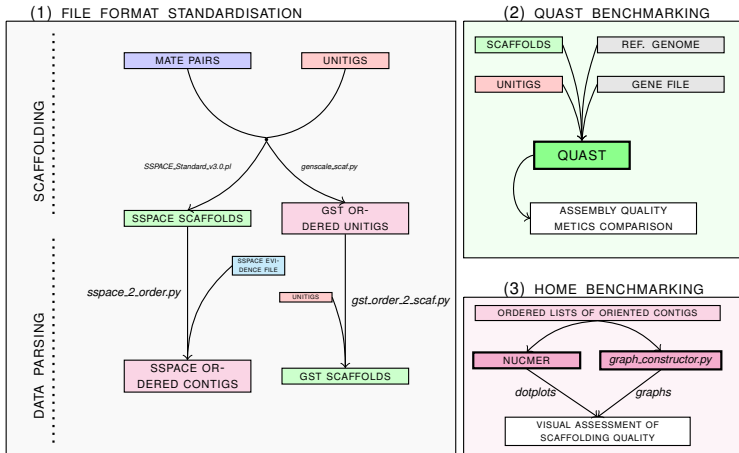
### Benchmarking workflow for the GST

### Results

Large repeats  
Short repeats

### Perspectives

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



# RESULTS

## 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

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 assemble because too many conflicting links exists 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

# LARGE REPEATS - CHLOROPLASTIC GENOMES

## 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

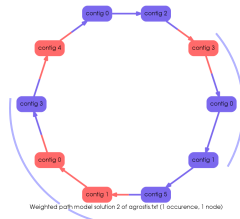
EXCELLENT RESULTS ARE OBTAINED FOR DATA SETS WITH LARGE REPEATS AND A SMALL NUMBER OF UNITIGS.



Golden standard / mapping solution of agrotis.tet (1 occurrence, 1 node)



SPRACE scaffolding solution of agrotis.tet



Weighted path model solution 2 of agrotis.tet (1 occurrence, 1 node)



Weighted path model solution 2 of agrotis.tet (1 occurrence, 1 node)

# SHORT REPEATS - BACTERIAL GENOMES

## SCAFFOLDING

ALEXANDRINA  
BODRUG

**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.

## Context

Some definitions  
Order and orient  
Challenging problem

## GST features

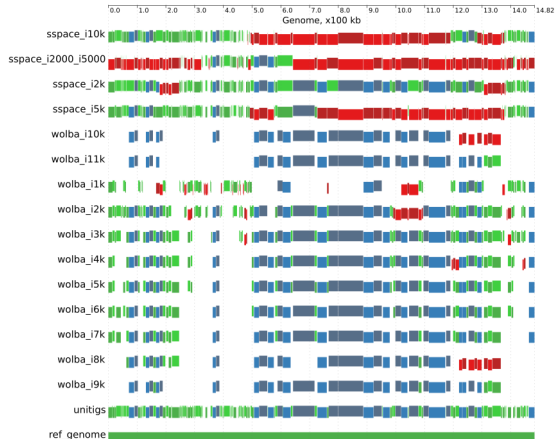
Scripting for the GST  
GST modeled graph  
Expected solution  
Scaffolding solutions

## Benchmarking workflow for the GST

## Results

Large repeats  
Short repeats

## Perspectives





# PERSPECTIVES

## 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

## DEVELOP - TEST - BENCHMARK

- Find strategies which solve more challenging data  
→ 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

# THANK YOU FOR YOU ATTENTION

## 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

## SCAFFOLDING: SAFETY COMES FIRST

