

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 23, 2015

Overview

SCAFFOLDING
BENCHMARK-
ING

ALEXANDRINA
BODRUG

Context

Some definitions
Order and orient

Genscale
scaffolding
tools features

The raw input data
GST modeled graph
Expected solution

Challenging
problem

Scripting for
the GST

Benchmarking
workflow for
tge GST

Results

Example of *Agrostis
stolonifera*

Perspectives

1 Context

- Some definitions
- Order and orient

2 Genscale scaffolding tools features

- The raw input data
- GST modeled graph
- Expected solution

3 Challenging problem

4 Scripting for the GST

5 Benchmarking workflow for tge GST

6 Results

- Example of *Agrostis stolonifera*

7 Perspectives

CONTEXT

SCAFFOLDING BENCHMARK- ING

ALEXANDRINA
BODRUG

Context

Some definitions
Order and orient

Genscale scaffolding tools features

The raw input data
GST modeled graph
Expected solution

Challenging problem

Scripting for the GST

Benchmarking workflow for the GST

Results

Example of *Agrostis
stolonifera*

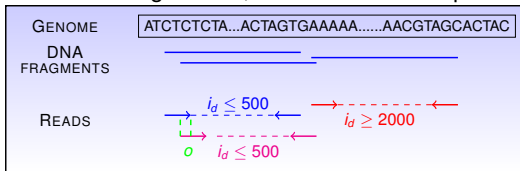
Perspectives

"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".

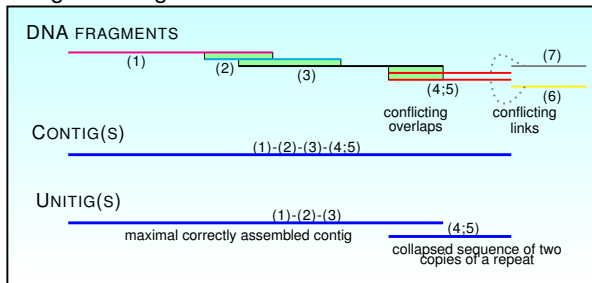
Hudson *et al.* 2002

SOME DEFINITIONS

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



...reads are assembled though high-confidence overlappings into contigs or unitigs.



ORDER AND ORIENT

SCAFFOLDING BENCHMARKING

ALEXANDRINA BODRUG

Context

Some definitions

Order and orient

Genscale scaffolding tools features

The raw input data

GST modeled graph

Expected solution

Challenging problem

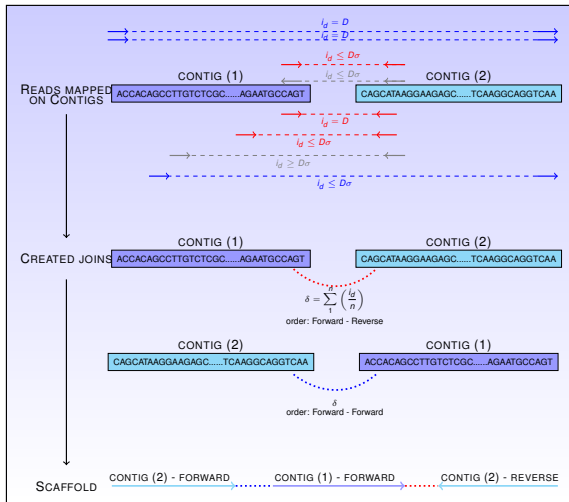
Scripting for the GST

Benchmarking workflow for tge GST

Results

Example of Agrostis
stolonifera

Perspectives



SCAFFOLDING BENCHMARK- ING

ALEXANDRINA
BODRUG

Context

Some definitions
Order and orient

Genscale scaffolding tools features

The raw input data
GST modeled graph
Expected solution

Challenging problem

Scripting for the GST

Benchmarking workflow for the GST

Results

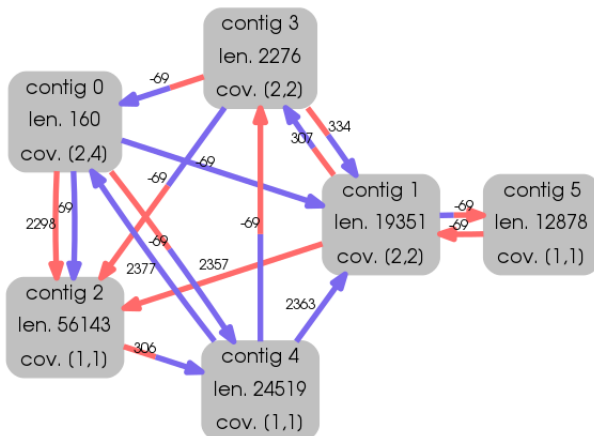
Example of Agrostis
stolonifera

Perspectives

- uses unitigs instead of contigs to better compute unitig coverage
- uses unitig coverages to duplicated regions
- several models exist, their common point is that for each unitig occurrence they create a node
- ... and for each unitig orientation, a different node is yet again created

THE RAW INPUT DATA

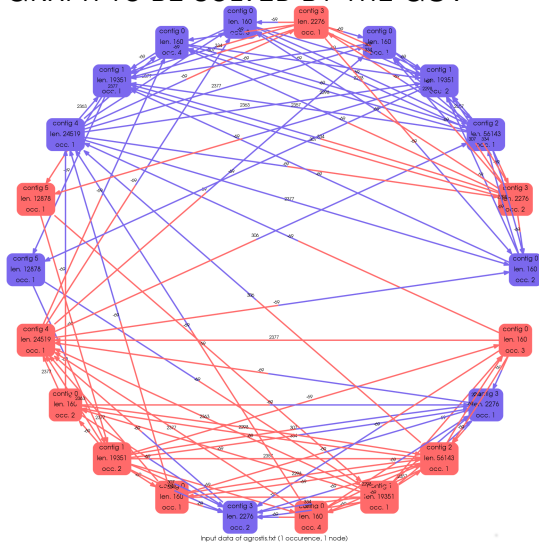
RAW INPUT DATA OF AGROSTIS STOLONIFERA



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

THE RAW INPUT DATA

GRAPH TO BE SOLVED BY THE GST



SCAFFOLDING BENCHMARK- ING

ALEXANDRINA BODRUG

Context

Some definitions
Order and orient

Genscale scaffolding tools features

The raw input data
GST modeled graph
Expected solution

Challenging problem

Scripting for the GST

Benchmarking workflow for tge GST

Results

Example of Agrostis
stolonifera

Perspectives

THE RAW INPUT DATA

SCAFFOLDING BENCHMARK- ING

ALEXANDRINA BODRUG

Context

Some definitions
Order and orient

Genscale scaffolding tools features

The raw input data
GST modeled graph
Expected solution

Challenging problem

Scripting for the GST

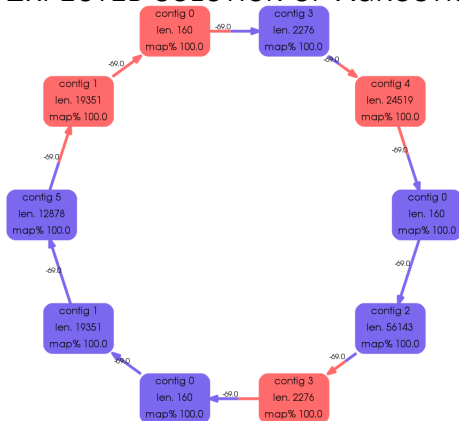
Benchmarking workflow for tge GST

Results

Example of Agrostis
stolonifera

Perspectives

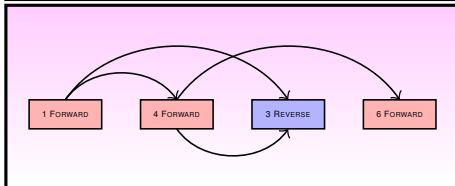
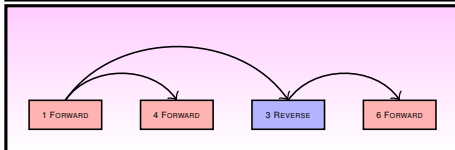
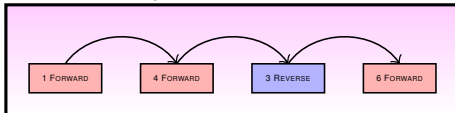
EXPECTED SOLUTION OF AGROSTIS STOLONIFERA



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

CHALLENGING PROBLEM

What would you do in these situations?



SCRIPTING FOR THE GST

SCAFFOLDING BENCHMARK- ING

ALEXANDRINA
BODRUG

Context

Some definitions
Order and orient

Genscale scaffolding tools features

The raw input data
GST modeled graph
Expected solution

Challenging problem

Scripting for the GST

Benchmarking workflow for the GST

Results

Example of *Agrostis
stolonifera*

Perspectives

- 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`

BENCHMARKING WORKFLOW FOR TGE GST

SCAFFOLDING BENCHMARK- ING

ALEXANDRINA
BODRUG

Context

Some definitions
Order and orient

Genscale scaffolding tools features

The raw input data
GST modeled graph
Expected solution

Challenging problem

Scripting for the GST

Benchmarking workflow for tge GST

Results

Example of Agrostis
stolonifera

Perspectives

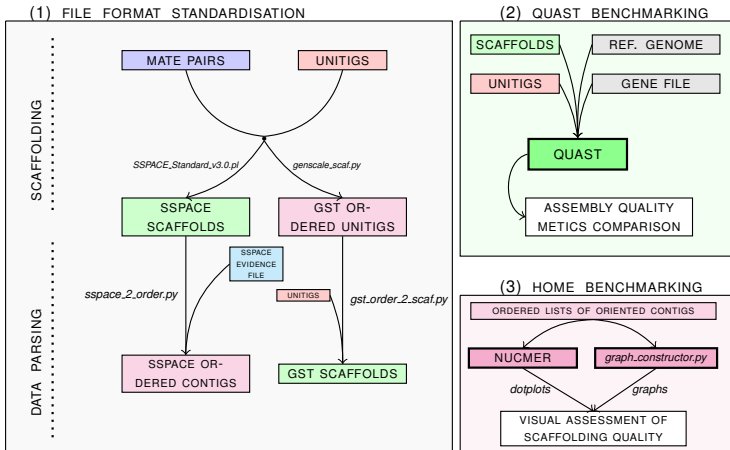


Figure : Benchmarking workflow

RESULTS

SCAFFOLDING BENCHMARK- ING

ALEXANDRINA
BODRUG

Context

Some definitions
Order and orient

Genscale scaffolding tools features

The raw input data
GST modeled graph
Expected solution

Challenging problem

Scripting for the GST

Benchmarking workflow for the GST

Results

Example of *Agrostis
stolonifera*

Perspectives

- Genomes with big repeated regions were solved a lot better than 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

EXAMPLE OF AGROSTIS STOLONIFERA

SCAFFOLDING BENCHMARKING

ALEXANDRINA BODRUG

Context

Some definitions
Order and orient

Genscale scaffolding tools features

The raw input data
GST modeled graph
Expected solution

Challenging problem

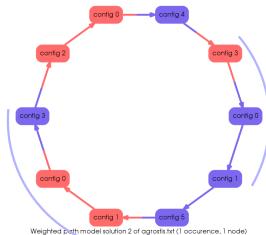
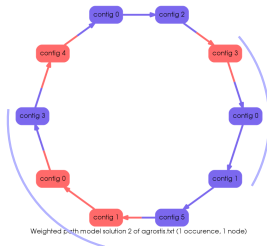
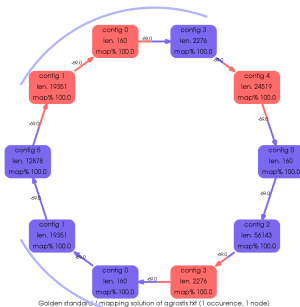
Scripting for the GST

Benchmarking workflow for tge GST

Results

Example of *Agrostis
stolonifera*

Perspectives



PERSPECTIVES

SCAFFOLDING BENCHMARK- ING

ALEXANDRINA
BODRUG

Context

Some definitions
Order and orient

Genscale scaffolding tools features

The raw input data
GST modeled graph
Expected solution

Challenging problem

Scripting for the GST

Benchmarking workflow for the GST

Results

Example of Agrostis
stolonifera

Perspectives

- Find strategies which solve more challenging data (flow model)
- Scaffold bacterial data
- Test the GST with real data

SCAFFOLDING BENCHMARK- ING

ALEXANDRINA
BODRUG

Context

Some definitions
Order and orient

Genscale scaffolding tools features

The raw input data
GST modeled graph
Expected solution

Challenging problem

Scripting for the GST

Benchmarking workflow for the GST

Results

Example of *Agrostis
stolonifera*

Perspectives

Thanks!

The End