

Reference Assisted Assembly and Annotation of the *Octopus vulgaris* Genome



Koen Herten^{1,2}, Gregory Maes^{1,2}, Eve Seuntjes³,
Fiorito Graziano⁴ and Joris R Vermeesch^{1,2}

1) Genomics Core, UZ Leuven, KU Leuven
2) Laboratory for Cytogenetics and Genome Research, Center for Human Genetics, KU Leuven
3) Animal Physiology and Neurobiology, KU Leuven
4) Stazione Zoologica Anton Dohrn

KU LEUVEN

Introduction

Compared to any other species, the *Octopus vulgaris* or common octopus appears to be utterly different, with its eight prehensile arms and its complex nervous system, enabling clever problem-solving and observational learning abilities and amazing physical behaviours such as millisecond colour and shape change (camouflage). Strikingly, the genome encoding these alien features turned out to be equally alien.

Materials & Methods



Octopus vulgaris
Estimated genome size: between 2.4-4.6GB
Estimated number of genes: 33,000



Sequencing:
213Gb (Giga Bases (A,C,G,T))
+90x coverage
HiSeq1500 PE 95bp
Insert sizes: 170bp, 250bp, 500bp, 800bp



Computing:
Thinking: Ivybridge, 2x10 cores per node, 124GB RAM per node
Cerebro: Ivybridge, 10 cores per node, 250GB RAM shared



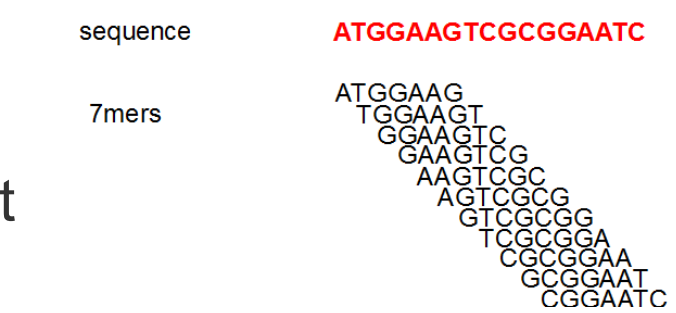
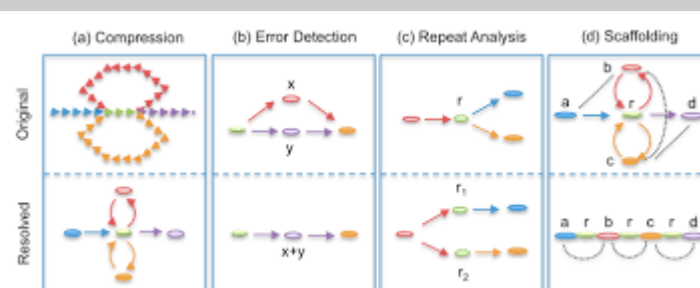
Can the genome be assembled by using the reference of the *Octopus bimaculoides*, and computationally annotated using the genes of the *Octopus bimaculoides*?

Workflow

Requirements*

Explanation

Results*

Pre-computation	Tool: SGA 0.10.14 cput: 55h + 10h mem: 105GB + 64GB Thinking, 1 node Input: raw (txt, gz) 165GB Output: txt 20GB	<ul style="list-style-type: none">Estimation of genome sizek-mer analysisEstimation of GC contentHeterozygosity estimation 	<ul style="list-style-type: none">Estimated genome size: 2.3GBEstimated GC content: 36%Estimated Heterozygosity Level: 1/30Best k-mer: 81														
De novo Assembly	Tool: AbySS 1.9.0 cput: 2,638h mem: 351GB Cerebro, 5 nodes Input: raw (txt, gz) 165GB Output: txt 3.4GB, total: 53GB	<ul style="list-style-type: none">De Bruijn graph construction of raw dataCreation of ContigsCreation of Scaffolds 	<table><tr><td>Estimated Genome Size</td><td>3.39GB</td></tr><tr><td># Scaffolds</td><td>7,452,219</td></tr><tr><td>N50</td><td>3072bp</td></tr><tr><td>L50</td><td>214,551</td></tr><tr><td>N75</td><td>1658bp</td></tr><tr><td>L75</td><td>479,686</td></tr><tr><td>#N's/100kbp</td><td>2551.2</td></tr></table>	Estimated Genome Size	3.39GB	# Scaffolds	7,452,219	N50	3072bp	L50	214,551	N75	1658bp	L75	479,686	#N's/100kbp	2551.2
Estimated Genome Size	3.39GB																
# Scaffolds	7,452,219																
N50	3072bp																
L50	214,551																
N75	1658bp																
L75	479,686																
#N's/100kbp	2551.2																
Reduction & Scaffolding	Tool: Redundans 0.12 cput: 470h mem: 293GB Cerebro, 5 nodes Input: raw (txt, gz) 165GB + txt 3.4GB Output: txt 2.3GB, total: 449GB	<ul style="list-style-type: none">Removal of high similar Contigs/ScaffoldsRe-Scaffolding using reduced Contig set and raw data	<table><tr><td>Estimated Genome Size</td><td>2.36GB</td></tr><tr><td># Scaffolds</td><td>1,290,144</td></tr><tr><td>N50</td><td>3524bp</td></tr><tr><td>L50</td><td>179,163</td></tr><tr><td>N75</td><td>1931bp</td></tr><tr><td>L75</td><td>394,315</td></tr><tr><td>#N's/100kbp</td><td>1126.2</td></tr></table>	Estimated Genome Size	2.36GB	# Scaffolds	1,290,144	N50	3524bp	L50	179,163	N75	1931bp	L75	394,315	#N's/100kbp	1126.2
Estimated Genome Size	2.36GB																
# Scaffolds	1,290,144																
N50	3524bp																
L50	179,163																
N75	1931bp																
L75	394,315																
#N's/100kbp	1126.2																
Reference Assisted Scaffolding	Tool: Blast 2.5.0 + Chromosomer 0.1.3 cput: 51m + 30m mem: 16GB + 2GB Thinking, 1 node Input: txt 3.4GB + txt 2.3GB Output: txt 1.7GB	<ul style="list-style-type: none">Find location of scaffolds on the genome of the <i>Octopus bimaculoides</i>Connect scaffold-set to chromosome level scaffolds	<table><tr><td>Estimated Genome Size</td><td>1.85GB</td></tr><tr><td># Scaffolds</td><td>84,018</td></tr><tr><td>N50</td><td>276,442bp</td></tr><tr><td>L50</td><td>1,602</td></tr><tr><td>N75</td><td>61,698bp</td></tr><tr><td>L75</td><td>4,939</td></tr><tr><td>#N's/100kbp</td><td>18,410.79</td></tr></table>	Estimated Genome Size	1.85GB	# Scaffolds	84,018	N50	276,442bp	L50	1,602	N75	61,698bp	L75	4,939	#N's/100kbp	18,410.79
Estimated Genome Size	1.85GB																
# Scaffolds	84,018																
N50	276,442bp																
L50	1,602																
N75	61,698bp																
L75	4,939																
#N's/100kbp	18,410.79																
Repeat Masking	Tool: WindowMasker 2.2.22 cput: 50m mem: 2GB Thinking, 1 node Input: txt 1.7GB Output: txt 1.7GB	<ul style="list-style-type: none">Find an mask repeat regions (repeats of small fragments of nucleotides)	<ul style="list-style-type: none">No stats of the repeat regions are currently available														
Annotation	Tool: splign 2.0.8 cput: 50m Mem: 2.2GB Thinking, 1 node Input: txt 1.7 GB + txt 65MB Output: 27MB	<ul style="list-style-type: none">Mapping of genes of the <i>Octopus bimaculoides</i> to the constructed <i>Octopus vulgaris</i> genomeFinding the locations of the genes in the <i>Octopus vulgaris</i> genome	<ul style="list-style-type: none">Found 15,036 possible gene locations15,036 of 38,556 <i>Octopus bimaculoides</i> genes have a location with high similarity on the new constructed <i>Octopus vulgaris</i> genome.														

* All measurements and results are obtained by using the standard parameters for all tools.

Discussion

- High Performance Computing is needed in order to solve the Big Data problem of *de novo* genome assembly. The Big Data in this projects measures around 1TB including raw data, intermediate files and the end results. Computing time exceeded over 3,000h, the needed memory over 350GB.
- The genome of the *Octopus vulgaris* appears to be alien, due to its high heterozygosity level. This biological problem introduces a big computational challenge, therefor this non-standard assembly workflow was used.
- The use of an evolutionary close species (*Octopus bimaculoides*) for a Reference Assisted Scaffolding, resulted in a genome with less scaffolds. Due to the evolutionary distance and the quality of the *Octopus bimaculoides* genome, this step could only partly solve this complex problem.
- Computational Annotation of the genome, by using the genes of an evolutionary close species (*Octopus bimaculoides*), resulted in a big set of high similar genes between the species. Expected is that a high number of genes are not found due to evolutionary distance and possible novel genes.
- The quality of the assembled genome can be improved by further parameter optimization, or by the generation of new data with the use of novel techniques like PacBio.

Acknowledgement

The computational resources and services used in this work were provided by the VSC (Flemish Supercomputer Center), funded by the Research Foundation Flanders (FWO) and the Flemish Government – department EWI