

# NOVOPlasty: De novo assembly of organelle genomes from whole genome data

Nicolas Dierckxsens<sup>1\*</sup>, Patrick Mardulyn<sup>1,2</sup> and Guillaume Smits<sup>1,2,3</sup>

<sup>1</sup>Interuniversity Institute of Bioinformatics in Brussels (IB2), Université Libre de Bruxelles and Vrije Universiteit Brussel, Triomaan CP 263, 1050 Brussels, Belgium,

<sup>2</sup>Evolutionary Biology and Ecology Unit, CP 160/12, Faculté des Sciences, Université Libre de Bruxelles, Av. F. D. Roosevelt 50, B-1050 Brussels, Belgium, <sup>3</sup>Genetics, Hôpital Universitaire des Enfants Reine Fabiola, Université Libre de Bruxelles, Brussels, Belgium and <sup>4</sup>Center for Medical Genetics, Hôpital Erasme, Université Libre de Bruxelles, Route de Lennik 808, 1070 Brussels, Belgium

\*To whom correspondence should be addressed. Tel: +32 0472 986806; Email: [nicolasdierckxsens@hotmail.com](mailto:nicolasdierckxsens@hotmail.com)

## Reference validation

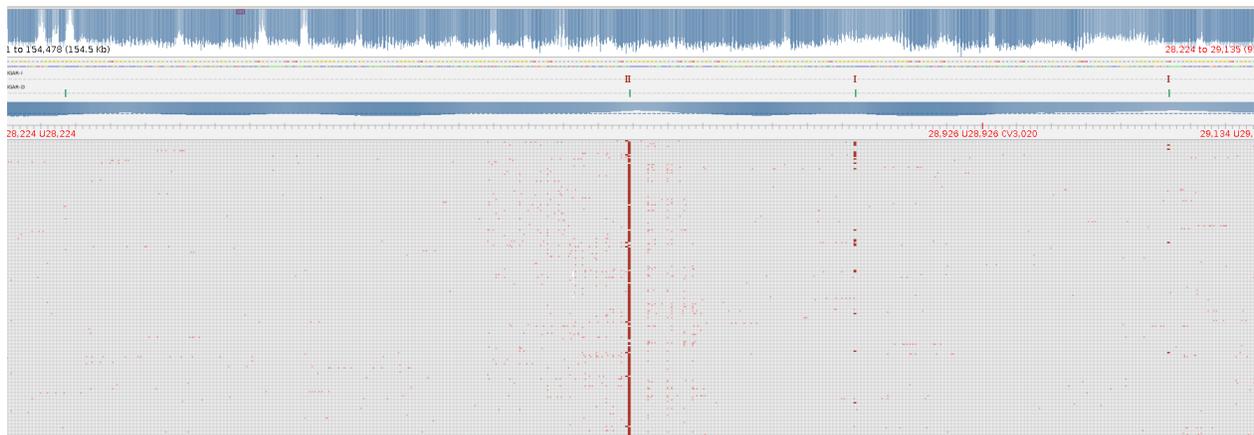
An accurate reference genome is essential to validate the quality of the assemblies. For the organelle genomes of *Oryza sativa*, *Arabidopsis thaliana* and *Homo sapiens* there are reference genomes available in GenBank (1). Nevertheless these reference genomes will not always be perfect because of small variances between individuals. Therefore each reference genome had to be validated manually by aligning (using bowtie2 (2)) all the reads of each dataset to the respective reference genome. Each variance was visually detected using Tablet (3) and corrected in the GenBank reference.

### *Homo sapiens*

No variance was found between GenBank entry X93334.1 and our dataset.

### *Arabidopsis thaliana*

The same GenBank entry AP000423.1 was used for both datasets SRR1174256 and SRR1810277. Two deletions of 1 bp were detected for SRR1174256 and one insertion of 1 bp for SRR1810277. All three were variations in the length of Single Nucleotide Repeats (SNR). The visual detection of the insertion is shown in Figure S1.



**Figure S1** | Visual detection of a single nucleotide insertion in a T homopolymer between dataset SRR1810277 and GenBank entry AP000423.1

*Oryza sativa*

No variance was found between GenBank entry KM103369.1 and dataset SRR1328237. Between GenBank entry KM088022.1 and dataset ERR477442, there was a difference of 19 nucleotides, consisting out of 10 Single-Nucleotide Polymorphisms (SNP's), 3 single bp insertions, a 4 bp long deletion and 2 single bp deletions. Figure S2 show an example of a SNP and Figure S3 shows the 4 bp long deletion.



**Figure S2 |** Visual detection of a SNP (A to C) between dataset ERR477442 and GenBank entry KM088022.1



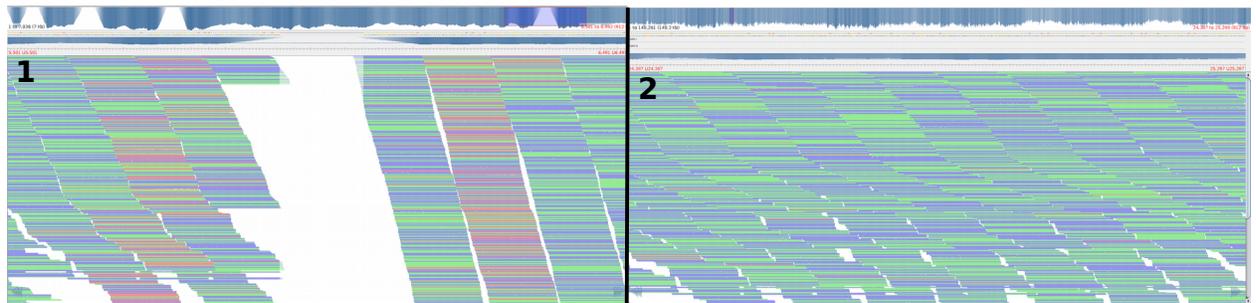
**Figure S3 |** Visual detection of a 4 bp deletion (AGGC) between dataset ERR477442 and GenBank entry KM088022.1

## Avicennia marina assembly comparison

There is no reference available for *Avicennia marina*, therefore we visually compared the differences between the CLC and the NOVOPlasty assemblies. This was done by realigning (with bowtie2) all the reads from the dataset to both of the assembly results and visually inspect the alignment with Tablet. The CLC assembly contained seven SNP's and was missing 1043 bp in comparison to the NOVOPlasty assembly. All missing regions and the 7 SNP's were confirmed by inspecting the alignment.



**Figure S4** | Visual confirmation of 4 SNP's between the CLC assembly and the *Avicennia marina* dataset.



**Figure S5** | Visual confirmation of 27 bp gap in the CLC assembly. Reads colored in blue or green are successfully paired together, reads in red are unpaired. (1) CLC assembly. (2) NOVOPlasty assembly.

## Benchmarking results

The benchmarking study consists out of 2 mitochondrial genomes and 5 chloroplast genomes. The complete set of results can be found in the following 7 tables. Table S8 contains the results of the *G. intermedia* dataset with a read length of 126 bp and an average coverage depth of 301. These results were excluded from the benchmarking study, due to lack of memory capacity for the MIRA and MITObim assemblies.

**Table S1 |** Benchmarking results for the assembly of the *Avicennia marina* chloroplast.

		<b><i>Avicennia marina</i> chloroplast</b>				
		NOVOPlasty	MIRA	MITObim	SOAPdenovo2	CLC
Duration	(min)	<b>6</b>	169	549	10	11
System+user time	(min)	<b>6</b>	319	1792	24	111
Memory	(GB)	7,4	21,7	5,2	17	<b>2,3</b>
Disk space	(GB)	<b>0,1</b>	30	24,5	<b>0,1</b>	2,2
Total contigs		<b>1</b>	172	87	47 449	249 654
Chloroplast contigs(scaffolds)		<b>1</b>	40	1	95	22(8)
Genome coverage		<b>100%</b>	65%	99,63%	86%	99,3%
Mismatches	(bp)	<b>0</b>	462		220	7
Ambiguous nucleotides (bp)		<b>4</b>	198		<b>0</b>	<b>0</b>
Accuracy		<b>100%</b>	99,52%	99,98%	99,83%	99,99%

**Table S2 |** Benchmarking results for the assembly of the *Oryza sativa* chloroplast (dataset SRR1328237).

		<b><i>Oryza sativa</i> (SRR1328237) chloroplast</b>				
		NOVOPlasty	MIRA	MITObim	SOAPdenovo2	CLC
Duration	(min)	<b>10</b>	89	18	84	21
System+user time	(min)	<b>10</b>	212	36	102	71
Memory	(GB)	4,4	12	<b>1,3</b>	14	1,5
Disk space	(GB)	<b>0,1</b>	27	3,3	0,9	1,7
Total contigs		<b>1</b>	525	3	261 800	186 057
Chloroplast contigs(scaffolds)		<b>1</b>	54	3	203	3(7)
Genome coverage		<b>100%</b>	31,82%	8,30%	39%	99,96%
Mismatches	(bp)	<b>0</b>	136	2	222	34
Ambiguous nucleotides (bp)		<b>1</b>	231	234	<b>0</b>	<b>0</b>
Accuracy		<b>100%</b>	99,68%	99,98%	99,58%	99,97%

**Table S3 |** Benchmarking results for the assembly of the *Oryza sativa* chloroplast (dataset ERR477442).

		<b><i>Oryza sativa</i> (ERR477442) chloroplast</b>				
		NOVOPlasty	MIRA	MITObim	SOAPdenovo2	CLC
Duration	(min)	5	61	6	<b>4</b>	9
System+user time	(min)	5	115	<b>7</b>	11	52
Memory	(GB)	4,3	11	<b>0,8</b>	11	1,3
Disk space	(GB)	<b>0,1</b>	20,5	0,9	0,2	0,9
Total contigs		<b>1</b>	1080	9	13 263	80 695
Chloroplast contigs(scaffolds)		<b>1</b>	35	9	77	7(3)
Genome coverage		<b>100%</b>	32,88%	4,82%	85%	99,97%
Mismatches	(bp)	<b>0</b>	37	142	27	2
Ambiguous nucleotides (bp)		<b>0</b>	157	154	<b>0</b>	<b>0</b>
Accuracy		<b>100%</b>	99,92%	97,90%	99,98%	<b>100%</b>

**Table S4** | Benchmarking results for the assembly of the *Arabidopsis thaliana* chloroplast (dataset SRR1174256).

<b><i>Arabidopsis thaliana</i> (SRR1174256) chloroplast</b>						
		NOVOPlasty	MIRA	MITObim	SOAPdenovo2	CLC
Duration	(min)	14	217	1389	43	<b>6</b>
System+user time	(min)	14	367	2001	56	60
Memory	(GB)	7,6	28	14	9	<b>1,3</b>
Disk space	(GB)	<b>0,1</b>	38	75	0,6	2,5
Total contigs		<b>1</b>	65 220	1	35 686	26 104
Chloroplast contigs(scaffolds)		<b>1</b>	41	0	200	8 (3)
Genome coverage		<b>100%</b>	33,05%	0%	56,19%	99,98%
Mismatches	(bp)	1	1022	/	<b>0</b>	4
Ambiguous nucleotides (bp)		1	367	/	<b>0</b>	<b>0</b>
Accuracy		<b>100%</b>	98%	/	<b>100%</b>	<b>100%</b>

**Table S5** | Benchmarking results for the assembly of the *Arabidopsis thaliana* chloroplast (dataset SRR1810277).

<b><i>Arabidopsis thaliana</i> (SRR1810277) chloroplast</b>						
		NOVOPlasty	MIRA	MITObim	SOAPdenovo2	CLC
Duration	(min)	9	220	3224	21	<b>6</b>
System+user time	(min)	9	398	4977	34	55
Memory	(GB)	6,9	27,5	23,5	10	<b>1,4</b>
Disk space	(GB)	<b>0,1</b>	123	126	0,6	1,8
Total contigs		<b>1</b>	72 115	398	13 614	97 204
Chloroplast contigs(scaffolds)		<b>1</b>	34	70	265	8 (3)
Genome coverage		<b>100%</b>	18,83%	68%	94%	99,98%
Mismatches	(bp)	1	119	43	232	10
Ambiguous nucleotides (bp)		<b>0</b>	19	65	<b>0</b>	<b>0</b>
Accuracy		<b>100%</b>	99,59%	99,96%	99,84%	99,99%

**Table S6** | Benchmarking results for the assembly of the *Homo sapiens* mitochondrion.

<b><i>Homo sapiens</i> mitochondrion</b>							
		NOVOPlasty	MIRA	MITObim	SOAPdenovo2	CLC	ARC
Duration	(min)	<b>4</b>	152	51	25	38	21
System+user time	(min)	<b>4</b>	745	51	69	269	33
Memory	(GB)	10	34,1	1,5	51	4,5	<b>0,6</b>
Disk space	(GB)	<b>0,1</b>	71	0,2	0,6	2,6	7
Total contigs		<b>1</b>	1161	<b>1</b>	4415	273 126	16
Chloroplast contigs(scaffolds)		<b>1</b>	<b>1</b>	<b>1</b>	2	<b>1</b>	16
Genome coverage		<b>100%</b>	99,98%	<b>100%</b>	<b>100%</b>	<b>100%</b>	42%
Mismatches	(bp)	<b>0</b>	1	<b>0</b>	<b>0</b>	<b>0</b>	/
Ambiguous nucleotides (bp)		<b>0</b>	4	<b>0</b>	<b>0</b>	<b>0</b>	/
Accuracy		<b>100%</b>	99,99%	<b>100%</b>	<b>100%</b>	<b>100%</b>	/

**Table S7 |** Benchmarking results for the assembly of the *Gonioctena intermedia* mitochondrion.

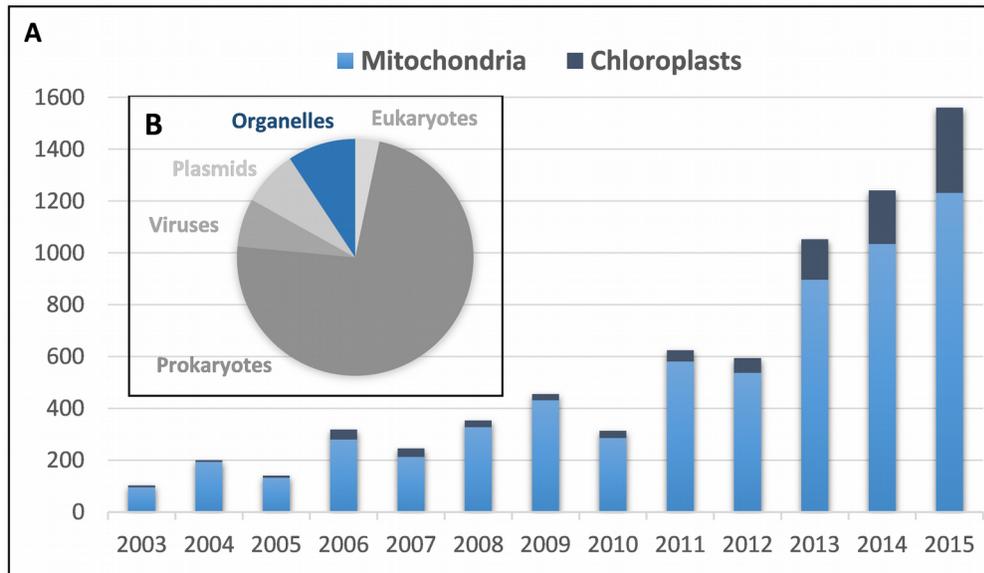
		<b><i>Gonioctena intermedia</i> mitochondrion</b>					
		NOVOPlasty	MIRA	MITObim	SOAPdenovo2	CLC	ARC
Duration	(min)	<b>11</b>	536	4777	19	51	586
System+user time	(min)	<b>11</b>	1451	8077	55	347	382
Memory	(GB)	15	57,6	63,4	27	5,1	<b>1,9</b>
Disk space	(GB)	<b>0,1</b>	144	418	0,9	3	12
Total contigs		<b>2</b>	3434	2221	3199	173 117	2502
Mitochondrial contigs		<b>1</b>	<b>1</b>	2	14	<b>1</b>	48
Genome coverage (complete mt)		92,75%	<b>93,66%</b>	93,29%	75,25%	89,96%	85,39%
Mismatches	(bp)	<b>0</b>	39	74	3	<b>0</b>	2
Ambiguous nucleotides (bp)		5	194	197	<b>0</b>	<b>0</b>	<b>0</b>
Accuracy		<b>100%</b>	99,77%	99,56%	99,98%	<b>100%</b>	99,99%
Genome coverage (minus tandem repeat)		<b>99,98%</b>	98,97%	98,95%	83,57%	99,85%	94,83%
Mismatches	(bp)	<b>0</b>	<b>0</b>	12	3	<b>0</b>	2
Ambiguous nucleotides (bp)		3	170	173	<b>0</b>	<b>0</b>	<b>0</b>
Accuracy		<b>100%</b>	<b>100%</b>	99,93%	99,97%	<b>100%</b>	99,99%

**Table S8 |** Assembly results for the mitochondrion of a second dataset of *Gonioctena intermedia* (excluded from the benchmark).

		<b><i>Gonioctena intermedia</i> mitochondrion 126 bp</b>					
		NOVOPlasty	MIRA	MITObim	SOAPdenovo2	SOAPdenovo2 (coverage depth -50%)	CLC
Duration	(min)	<b>14</b>			75	26	213
System+user time	(min)	<b>14</b>			162	68	1114
Memory	(GB)	27			54	29	<b>10,1</b>
Disk space	(GB)	<b>0,1</b>			5	1,8	5,3
Total contigs		<b>1</b>			1 796 373	552 116	590 150
Mitochondrial contigs		<b>1</b>	OUT	OUT	14	20	8(3)
Genome coverage (complete mt)		<b>93,32%</b>	OF	OF	27,74%	46,80%	88,70%
Mismatches	(bp)	<b>0</b>	MEMORY	MEMORY	14	22	13
Ambiguous nucleotides (bp)		3			<b>0</b>	<b>0</b>	2
Accuracy		<b>100%</b>			99,74%	99,74%	99,92%
Genome coverage (minus tandem repeat)		<b>99,99%</b>			30,80%	51,97%	99,85%
Mismatches	(bp)	<b>0</b>			14	22	13
Ambiguous nucleotides (bp)		1			<b>0</b>	<b>0</b>	2
Accuracy		<b>100%</b>			99,74%	99,74%	99,92%

## Extra Figures

### Annual deposition of organelle genomes in GenBank



**Figure S6 |** Organelle genomes in GenBank. (A) Total number of deposited genomes (15704) in GenBank as of 14 March 2016. (B) Annual deposition of mitochondrial and chloroplast genomes in GenBank since 2003. Statistics from the National Center for Biotechnology Information Genome Resources (<https://www.ncbi.nlm.nih.gov/genome/browse/> (13 March 2016, date last accessed)).

## References

1. Benson, D.A., et al. (2013) GenBank. *Nucleic Acids Research*, 41, D36D42.
2. Langmead, B., Salzberg, S. (2012) Fast gapped-read alignment with Bowtie2. *Nature Methods*, 9, 357-359.
3. Milne, I., Stephen, G., Bayer, M., Cock, P.J.A., Pritchard, L., Cardle, L., Shaw, P.D. and Marshall, D. (2013) Using Tablet for visual exploration of second-generation sequencing data. *Briefings in Bioinformatics*, 14(2), 193-202.