

# 1 A draft genome of grass pea (*Lathyrus* 2 *sativus*), a resilient diploid legume

3 Authors: Peter M. F. Emmrich<sup>1,2</sup>, Abhimanyu Sarkar<sup>1</sup>, Isaac Njaci<sup>1,2</sup>, Gemy George Kaithakottil<sup>3</sup>, Noel  
4 Ellis<sup>1</sup>, Christopher Moore<sup>4</sup>, Anne Edwards<sup>1</sup>, Darren Heavens<sup>3</sup>, Darren Waite<sup>3</sup>, Jitender Cheema<sup>1</sup>,  
5 Martin Trick<sup>1</sup>, Jonathan Moore<sup>1</sup>, Anne Webb<sup>5</sup>, Rosa Caiazzo<sup>5</sup>, Jane Thomas<sup>5</sup>, Janet Higgins<sup>3</sup>, David  
6 Swarbreck<sup>1</sup>, Shiv Kumar<sup>6</sup>, Sagadevan Mundree<sup>7</sup>, Matt Loose<sup>4</sup>, Levi Yant<sup>4</sup>, Cathie Martin<sup>1</sup>, Trevor L.  
7 Wang<sup>1</sup>.

8 <sup>1</sup> John Innes Centre, Norwich Research Park, Colney Lane, Norwich, NR4 7UH, United Kingdom

9 <sup>2</sup> Biosciences eastern and central Africa International Livestock Research Institute Hub, ILRI campus,  
10 Naivasha Road, P.O. 30709, Nairobi 00100, Kenya

11 <sup>3</sup> Earlham Institute, Norwich Research Park, Colney Lane, Norwich, NR4 7UZ, United Kingdom

12 <sup>4</sup> University of Nottingham, University Park, Nottingham, NG7 2RD, United Kingdom

13 <sup>5</sup> National Institute of Agricultural Botany, Huntingdon Road, Cambridge, CB3 0LE, GB

14 <sup>6</sup> International Center for Agricultural Research in the Dry Areas, Avenue Hafiane Cherkaoui, Rabat,  
15 Morocco

16 <sup>7</sup> Queensland University of Technology, 2 George St, Brisbane City QLD 4000, Australia

17

## 18 Abstract

19 We have sequenced the genome of grass pea (*Lathyrus sativus*), a resilient diploid (2n=14) legume  
20 closely related to pea (*Pisum sativum*). We determined the genome size of the sequenced European  
21 accession (LS007) as 6.3 Gbp. We generated two assemblies of this genome, i) Elv1 using Illumina  
22 PCR-free paired-end sequencing and assembly followed by long-mate-pair scaffolding and ii) Rbp  
23 using Oxford Nanopore Technologies long-read sequencing and assembly followed by polishing with  
24 Illumina paired-end data. Elv1 has a total length of 8.12 Gbp (including 1.9 billion Ns) and scaffold  
25 N50 59,7 kbp. Annotation has identified 33,819 high confidence genes in the assembly. Rbp has a  
26 total length of 6.2 Gbp (with no Ns) and a contig N50 of 155.7 kbp. Gene space assessment using the  
27 eukaryote BUSCO database showed completeness scores of 82.8 % and 89.8%, respectively.

28 **Keywords:** grass pea, *Lathyrus*, legume, beta-ODAP, genome

## 29 Introduction

30 Grass pea (*Lathyrus sativus* L.), is a diploid legume ( $2n=14$ ) first domesticated in the Balkan peninsula  
31 (Kislev 1989). It is grown for its seeds as food and feed, as well as for fodder, mainly by farmers in  
32 the Indian subcontinent as well as in northern African countries such as Ethiopia (Campbell 1997;  
33 Kumar et al. 2011). It is self-fertile, though under field conditions it is often cross pollinated by  
34 insects such as bees. Grass pea shows remarkable tolerance to environmental stress, including both  
35 drought and waterlogging (Yadav et al. 2006; Campbell 1997) making it a vital source of food and  
36 feed in times of scarcity, demonstrated repeatedly during its 8000 years of cultivation (Campbell  
37 1997). It is highly efficient at fixing atmospheric nitrogen in soils using rhizobial symbionts, meaning  
38 that it has low input requirements for fertilizers (Jiao et al. 2011; Drouin, Prévost, and Antoun 2000).  
39 For many, grass pea represents an ‘insurance crop’ as it survives and produces a yield under  
40 conditions, such as drought or flooding, where most other crops fail (Girma, Tefera, and Dadi 2011;  
41 Zhelyazkova et al. 2016; Silvestre et al. 2014; Yang and Zhang 2005; Vaz Patto et al. 2006) and it is  
42 often grown by poor farmers using minimal inputs to ensure a supply of food. Its seeds are high in  
43 protein (up to 30% w/w in dry seed) (Emmrich 2017), and it has the potential to provide nutritional  
44 food security in some of the most resource scarce and least developed regions in the world in the  
45 face of climate change and increasing population pressure (Sarkar et al. 2019).

46 The chief drawback of grass pea as a food has been the presence of an anti-nutritional factor,  $\beta$ -N-  
47 oxalyl-L- $\alpha,\beta$ -diaminopropionic acid ( $\beta$ -ODAP) in most tissues of the plant, including seeds. In  
48 conjunction with severe malnutrition, prolonged consumption of  $\beta$ -ODAP causes a neurological  
49 disorder, neurolathyrism, which results in permanent paralysis of the lower limbs in humans (Dufour  
50 2011; Kusama-Eguchi et al. 2014). This has led to bans on grass pea cultivation and commerce (Cohn  
51 and Streifler 1983), leading to underinvestment in research in this promising legume. The main  
52 thrust of grass pea research has been the elimination of  $\beta$ -ODAP from grass pea. A number of low  $\beta$ -  
53 ODAP grass pea cultivars have been developed by plant breeders, such as BioL-212 (Ratan),  
54 Mahateora, LS 8246, Prateek and Ceora (Kumar et al. 2011; Sawant, Jayade, and Patil 2011;  
55 Chakrabarti, Santha, and Mehta 1999; Santha and Mehta 2001; Tsegaye, Tadesse, and Bayable 2005;  
56 Siddique, Hanbury, and Sarker 2006), but no zero  $\beta$ -ODAP grass pea cultivar has been produced to  
57 date. The biosynthesis of  $\beta$ -ODAP in grass pea is partially understood (Lambein et al. 1993; Kuo and  
58 Lambein 1991; Ikegami et al. 1999; Ikegami et al. 1993; Malathi, Padmanab, and Sarma 1970) but  
59 critical gaps remain in our knowledge of the genetics of  $\beta$ -ODAP biosynthesis in grass pea.

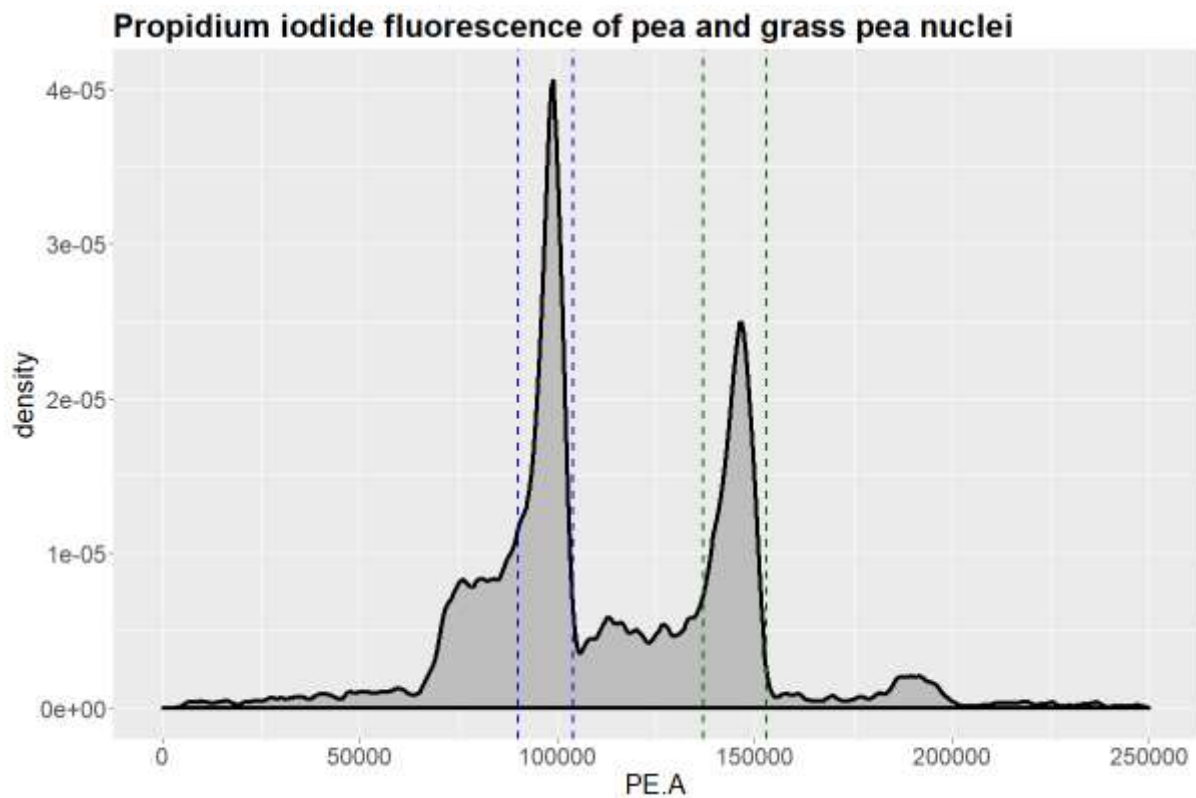
60 We present a draft genome sequence of a European accession (LS007) of grass pea which will be  
61 valuable for the identification of genes in the  $\beta$ -ODAP biosynthesis pathway, as well as in

62 identification and selection of traits for agronomic improvement. The data will allow comparative  
63 genomic analyses between legumes, help in the development of high quality genetic and physical  
64 maps for marker-assisted and genomic selection strategies and enable genome editing and TILLING  
65 platforms for grass pea improvement. The availability of this draft genome sequence will facilitate  
66 research on grass pea with the goal of developing varieties that fulfil its potential as a high protein,  
67 low input, resilient, climate smart crop, suitable for small-holder farmers.

## 68 Results

### 69 Genome size estimation

70 Using *Pisum sativum* as a standard with a genome size of 4.300 Gbp (Leitch et al. 2019), we  
71 undertook flow cytometry and estimated the genome size of grass pea genotype LS007 as 6.297 Gbp  
72  $\pm 0.039$  Gbp (stdev.). Peak coefficients of variance were below 3% in all three replicates (see Fig. 1).



73

74 *Figure 1. Results of a representative flow cytometry run, after gating to exclude cell debris events.*  
75 *Propidium iodide fluorescence amplitude (in arbitrary units) is shown against event density. The*  
76 *interval assumed to be pea nuclei is defined by blue dashed lines, the interval assumed to be grass*  
77 *pea nuclei is defined by green dashed lines. The experiment was replicated three times.*

78 We used this estimate to calculate genome sequencing depth of all our sequencing datasets. As this  
79 estimate is based on *Pisum sativum* as a standard, any inaccuracy in the *Pisum sativum* genome size  
80 estimate will affect the genome size estimate for LS007. The Kew Gardens c-value database records  
81 13 *Pisum sativum* samples with haploid genome sizes ranging from 3.724 Gbp to 5.782 Gbp (Leitch  
82 et al. 2019). Applied to the grass pea genome this would result in a genome size range of 5.456 Gbp  
83 to 8.471 Gbp.

## 84 Sequencing and Genome assembly

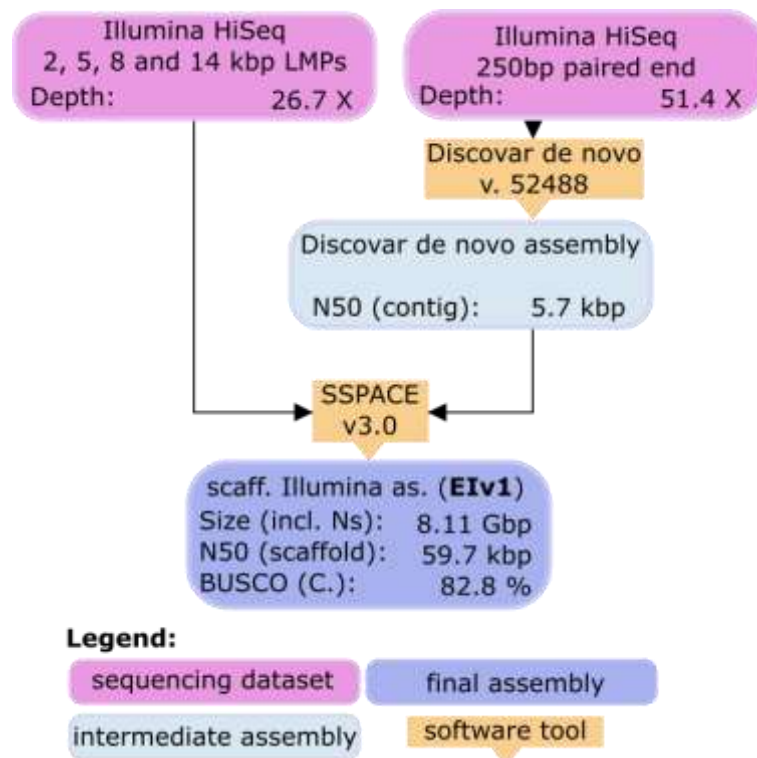
85 Figures 2 and 3 summarise the workflow that generated the datasets and the strategies used to  
86 assemble the two versions of the grass pea genome we discuss here, one based entirely on Illumina  
87 HiSeq paired-end data scaffolded using Illumina HiSeq long-mate-pair data (Elv1) and one based on  
88 Oxford Nanopore Technologies PromethION data assembled de novo, followed by polishing using  
89 the Illumina paired-end data.

### 90 Paired-end sequencing using Illumina platform

91 The paired end (PE) sequencing of LS007 genomic DNA was carried out using PCR-free libraries,  
92 followed by sequencing of the Long Mate Pair libraries (LMP) on the HiSeq Illumina platform (see  
93 Material and Methods for details).

### 94 Assembly

95 The final Elv1 assembly consisted of 669,893 contigs (minimum size 1 kbp), with a scaffold N50 value  
96 of 59,728 bases for a total assembly of 8.12 Gbp, including 1.9 Gbp of Ns in scaffolds.



97

98 *Fig. 2 Workflow of assembly and scaffolding of Elv1 from Illumina paired-end and LMP data*

## 99 Long-read sequencing using the PromethION platform

### 100 Sequence data

101 Sequence yields for each load of the flowcells are shown in Table 1. Subsequent loads on the same  
102 flowcell were separated by nuclease flushes. In total, 296.15 Gbp of sequence passed the quality  
103 filter, representing 47.01 X coverage of a 6.3 Gbp genome (see results section). Distributions of read  
104 lengths (post-filter) are shown below.

105

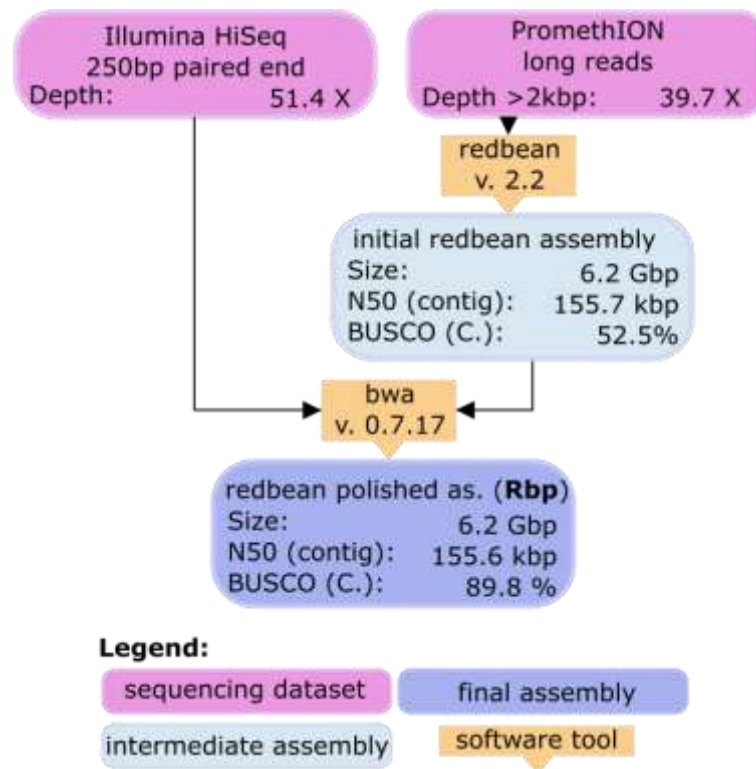
106 *Table 1. ONT PromethION sequencing yield*

Flowcell/Nuclease flush	DNA preparation	Library amount (ng)	Read N50 (kbp)	Yield (Gbp) (passed)
FC1 1st load	CN-AMP	250	7.84	51.84
FC1 2nd load	CN-NS-AMP	250	5.25	44.41
FC1 3rd load	QD-SRE	250	18.78	14.28
FC2	QD-SRE	250	22.84	57.48
FC3 1st load	CN-NS-AMP-SRE	270	22.31	41.7
FC3 2nd load	Mix of CN-NS-AMP-SRE & QD-SRE	324	22.81	23.04
FC4	CN-NS-AMP-SRE	400	29.08	32.7
FC5	CN-NS-AMP-SRE	400	24.71	30.7

107

### 108 Assembly

109 After filtering for reads >5 kbp, Redbean (Ruan and Li 2020) produced an assembly of 6.2 Gbp, based  
110 on 35.8 X coverage. The resulting assembly contained 162,985 contigs, with a contig N50 of 155,574  
111 bp. Following polishing with minimap2 (Li 2018) and bwa (Li 2013), the assembly was brought to a  
112 size of 6.237 Gbp in 162,994 contigs with an N50 of 157,998 bp.



113

114 *Figure 3. Workflow of assembly of Rbp from ONT data and polishing with Illumina paired-end data.*

## 115 Assembly comparison

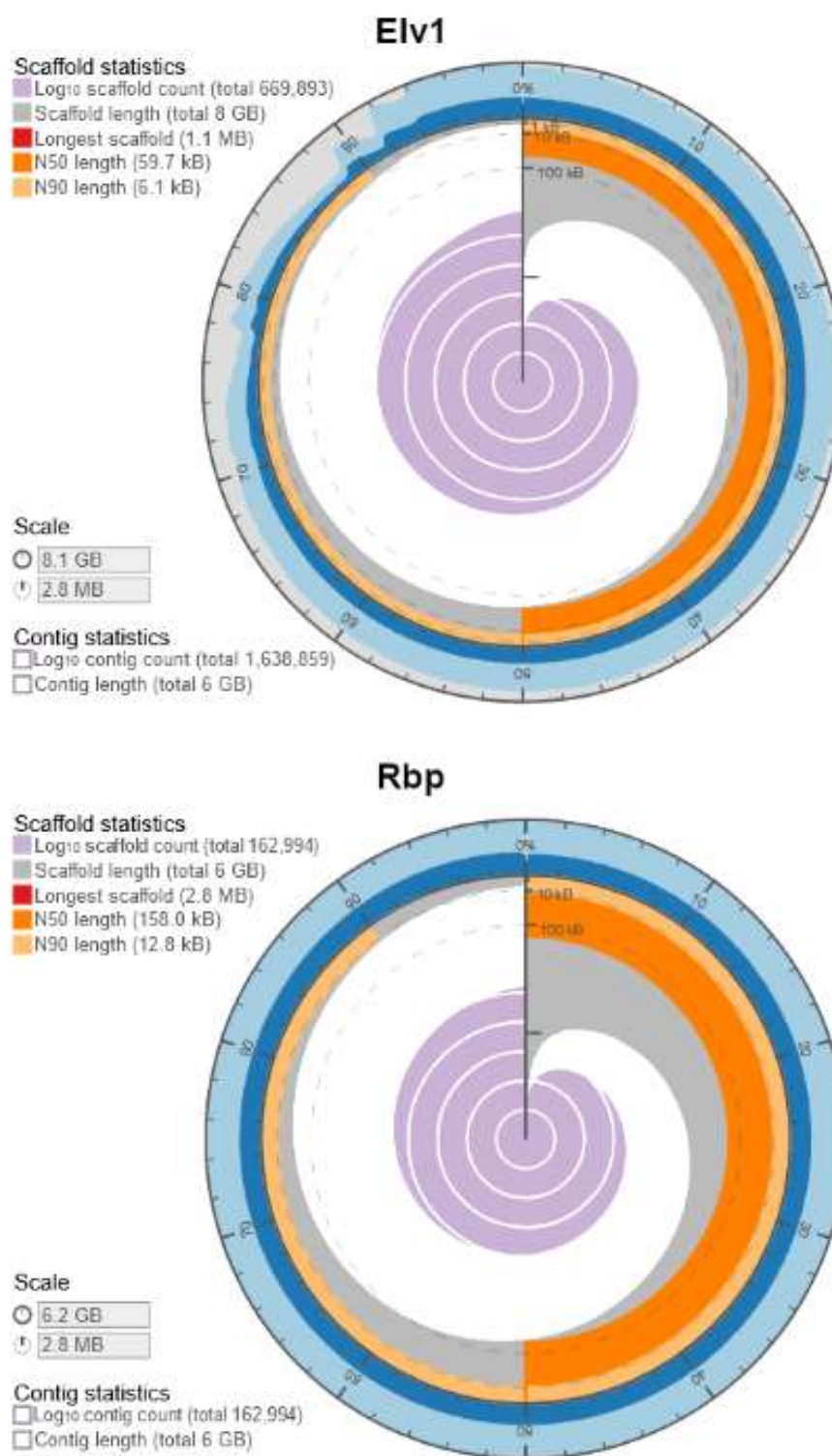
116 Comparison of the Elv1 and the polished Redbean assemblies revealed very similar number of ATGC  
 117 bases (approx. 6.2 Gbp), the size difference being mostly accounted for by the approximately 1.9  
 118 billion Ns in the Elv1 assembly (see Table 2). The GC fraction was 0.388 for the polished Redbean  
 119 assembly and 0.383 for the Elv1 assembly. Assembly statistics are shown in Figure 4.

120

121 *Table 2. Statistics of the Elv1 and Rbp assemblies.*

	<b>Elv1 assembly</b>	<b>polished Redbean assembly (Rbp)</b>
<b>A</b>	1,912,110,396	1,889,610,210
<b>T</b>	1,911,544,733	1,924,885,321
<b>G</b>	1,188,430,610	1,230,910,845
<b>C</b>	1,188,469,488	1,191,840,302
<b>N</b>	1,918,563,149	0
<b>ATGC bases</b>	6,200,555,227	6,237,246,678
<b>Total length</b>	8,119,118,376	6,237,246,678
<b>GC fraction</b>	38.3 %	38.8 %
<b>Longest scaffold/contig</b>	1,110,364	2,768,903
<b>N50</b>	59,728	157,998
<b>L50</b>	31,600	8,679
<b>No. of scaffolds/contigs</b>	669,893	162,994

122



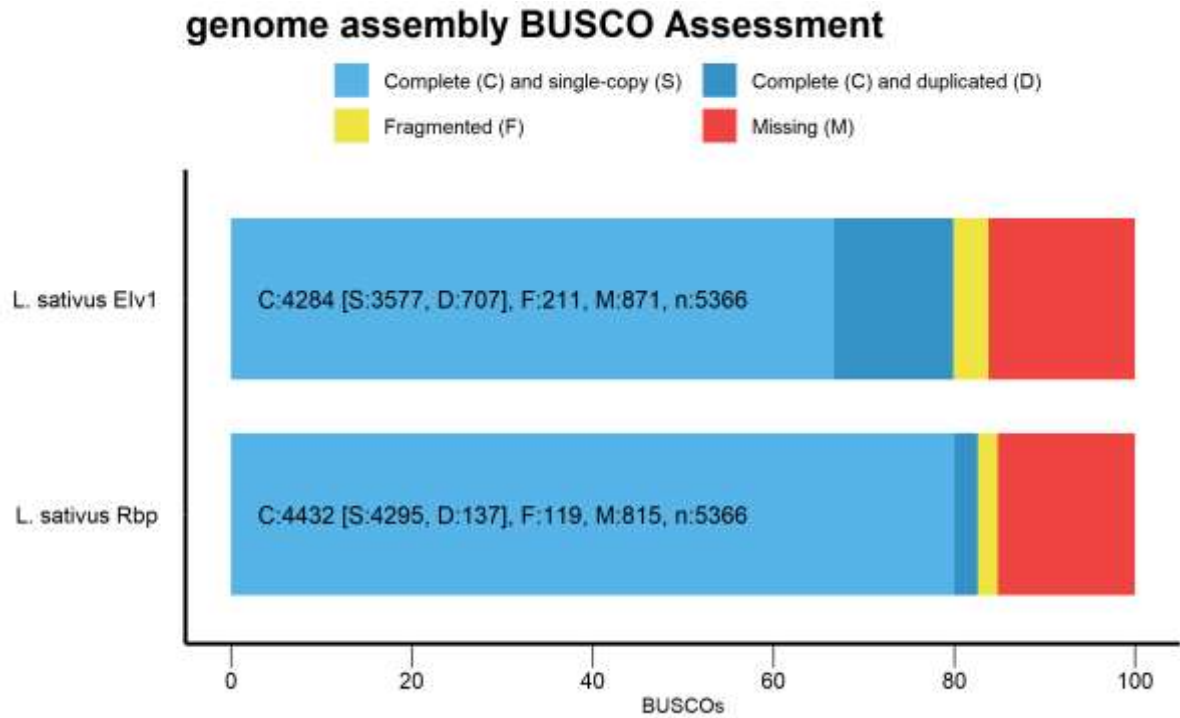
123

124 *Figure 4. Statistics of the Elv1 and Rbp assemblies of the LS007 genome, visualised using assembly-*  
125 *stats (Challis, 2015). Note the difference in scale for contig/scaffold size.*



126 BUSCO analysis

127 BUSCO (Benchmarking Universal Single-Copy Orthologs) analysis was carried out on both the  
 128 Illumina-only (Elv1) and long-read Redbean assembly polished with Illumina paired-end data (Rbp)  
 129 to assess the genome assemblies for gene space completeness (Simão et al. 2015). Results are  
 130 shown in Figure 5 and Table 3.



131

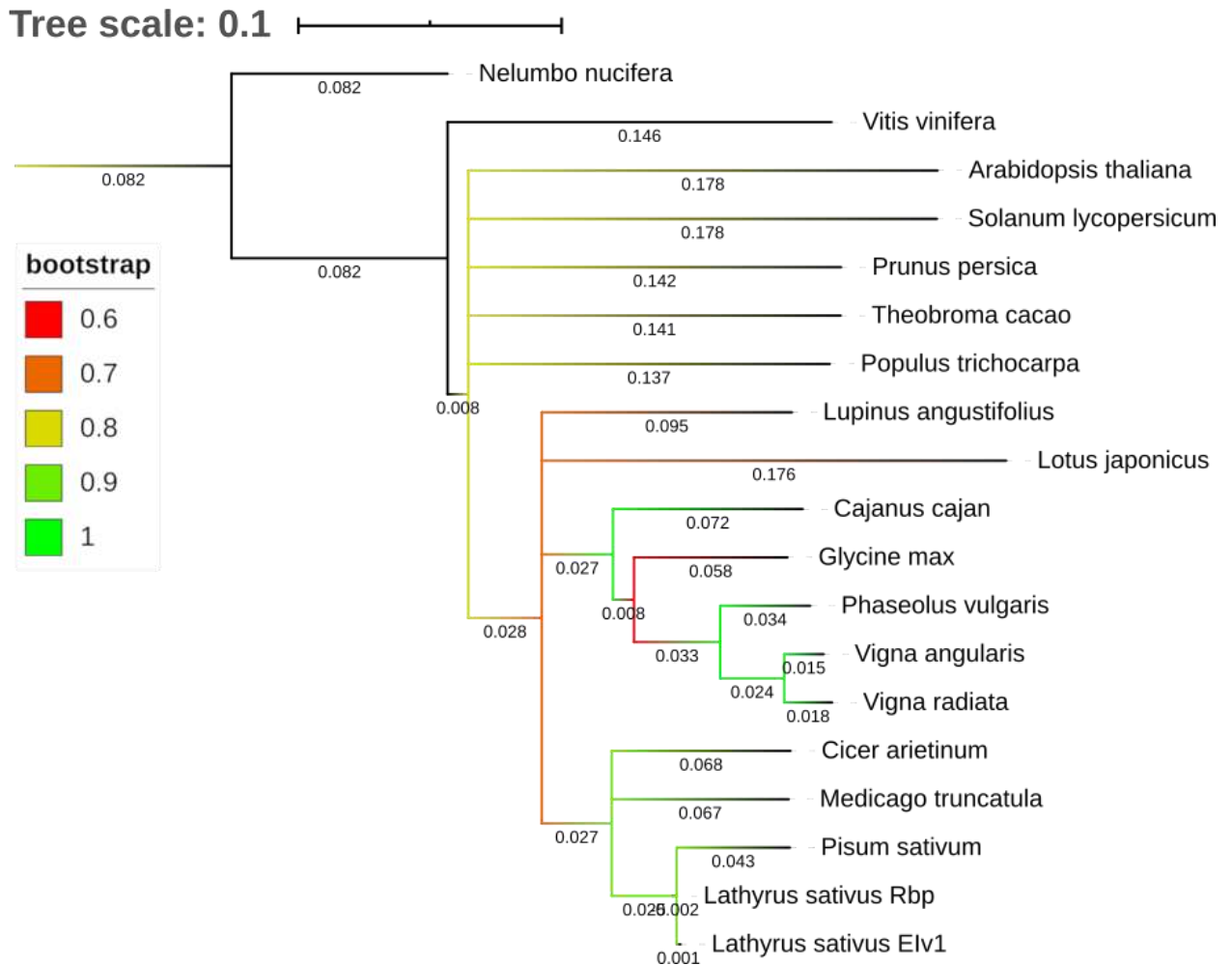
132 *Figure 5. BUSCO (v.4.0.4) assessment of the Lathyrus sativus LS007 genome Elv1 and Rbp assemblies*  
 133 *against the Fabales dataset. 5366 total BUSCO groups were searched.*

134

135 *Table 3. BUSCO analysis results using BUSCO:v4.0.4\_cv1 and odb10 databases*

Database	Fabales (n:5366)		Eudicots (n:2326)		Viridiplantae (n:425)		Eukaryote (n:255)	
	Elv1	Rbp	Elv1	Rbp	Elv1	Rbp	Elv1	Rbp
<b>Complete BUSCOs (C)</b>	<b>79.9%</b>	<b>82.6%</b>	<b>81.7%</b>	<b>85.7%</b>	<b>86.4%</b>	<b>88.4%</b>	<b>82.8%</b>	<b>89.8%</b>
- <b>Single copy BUSCOs (S)</b>	66.7%	80.0%	68.2%	82.7%	71.3%	84.9%	66.3%	82.0%
- <b>Duplicated BUSCOs (D)</b>	13.2%	2.6%	13.5%	3.0%	15.1%	3.5%	16.5%	7.8%
<b>Fragmented BUSCOs (F)</b>	<b>3.9%</b>	<b>2.2%</b>	<b>6.7%</b>	<b>3.0%</b>	<b>9.9%</b>	<b>3.8%</b>	<b>8.2%</b>	<b>3.9%</b>
<b>Missing BUSCOs (M)</b>	<b>16.2%</b>	<b>15.2%</b>	<b>11.6%</b>	<b>11.3%</b>	<b>3.7%</b>	<b>7.8%</b>	<b>9.0%</b>	<b>6.3%</b>

136 A phylogeny of grass pea (LS007) was determined in the context of 17 other plant species, based on  
137 the sequences of 10 BUSCO genes from the viridiplantae set (Figure 6). The sequences for these  
138 arbitrary genes are highly similar in both our assemblies causing them to be grouped together in the  
139 phylogeny.



140

141 *Figure 6. Phylogeny of both LS007 assemblies and 17 other plant species, based on 10 arbitrary*  
142 *BUSCO gene models from the viridiplantae set, re-rooted to Nelumbo nucifera. Colours represent*  
143 *branch bootstrap support. Plotted using iTOL (Letunic and Bork 2019).*

## 144 Repeat structure

145 We screened a subset (0.1 x coverage) of Illumina paired-end reads for repeated elements using the  
146 RepeatExplorer2 algorithm. Table 4 shows the portion of reads classified as repeats, along with  
147 literature values reported for *Lathyrus sativus* by Macas et al. (2015). As this analysis was done on  
148 the basis of the raw read data, it was independent of the assembly strategy used.

149

150 *Table 4. RepeatExplorer2 results summary.*

Repeat type					<i>L. sativus</i> LS007 (%)	<i>L. sativus</i> (Macas et al. 2015) (%)	
rDNA					0.59	1.69	
satellite					8.12	10.73	
Mobile element	Class I	LTR_unclassified			2.44	5.06	
		Max./SIRE			6.84	6.85	
		Ty_1/copia	Ale		0.11	0.02	
			Angela		0.26	0.2	
			Ivana		0.65	0.29	
			TAR		0.12	0.07	
			Tork		0.52	0.33	
		Ty_3/gypsy	Non- chrom.	OTA	Athila	3.87	3.11
					Tat (other)	0.18	0.51
					Ogre	37.32	45.46
		Chromovirus			4.55	3.33	
	Class II	Subclass I	TIR	EnSpm_CACTA	0.75	0.03	
				hAT	0.03	0.01	
				MuDR_Mutator	0.18	0.03	
		Subclass II	Helitron		0.11	0	
unclassified/no_evidence					4.13	2.64	
TRIM					0.00	0.03	
Tandem					0.00	0.86	
<b>total</b>					<b>70.78</b>	<b>81.25</b>	

151

152 **Annotation**

153 **Elv1 assembly**

154 We identified a total of 33,819 high confidence gene models in the Elv1 assembly of which 940 were  
 155 associated with repeats. On average, there were 1.05 transcripts associated with each gene, with a

156 mean cDNA size of 1,454.22 bp and mean transcript size (including introns) of 3,488.6 bp. There  
157 were, on average, 4.86 exons per transcript, with an average exon size of 299.17 bp.

158 A further 53,403 low confidence gene models were also identified, of which 9,577 were associated  
159 with repeats. On average, these had 1.03 transcripts associated with each gene, with a mean cDNA  
160 size of 752.91 bp and mean transcript size (including introns) of 2645.52 bp. They had, on average,  
161 2.47 exons per transcript, with an average exon size of 305.03 bp.

162 The number of genes within each category is shown in Table 5. The genome annotation summary  
163 statistics are given in Supplemental table S1.

164 *Table 5. Annotation biotype and gene confidence assignment for Elv1 assembly*

<b>Confidence level</b>	<b>Biotype</b>	<b>Gene count</b>
High	protein_coding	32,879
High	protein_coding_repeat_associated	940
Low	protein_coding	43,826
Low	protein_coding_repeat_associated	9,577

165

## 166 tRNA gene prediction

167 tRNA gene prediction was carried out using tRNAscan-SE (Lowe and Eddy 1997), and resulted in a  
168 total of 5765 tRNA genes predicted for the Elv1 assembly and 2801 tRNAs genes predicted for the  
169 polished Redbean (Rbp) assembly. Supplemental Table S2 reports the final number of coding  
170 transcripts per each rank.

171

## 172 Discussion

173 The grass pea genome follows the recent publication of the reference genome of *Pisum sativum*  
174 (Kreplak et al. 2019), a close relative of grass pea. Assembling the genome of grass pea is rendered  
175 more difficult due to its large size (we estimate 6.3 Gbp for cultivar LS007) and the comparative lack  
176 of genetic resources such as genetic and/or molecular maps that could be used for high-level  
177 scaffolding to the pseudochromosome level.

178 Our estimate for the genome size of grass pea (6.3 Gbp) measured by flow cytometry contrasts with  
179 the range of values quoted for *Lathyrus sativus* genome size in the literature, including the value  
180 listed in the Kew Plant C-Value database of 8.2 Gbp and the ranges given by Ghasem et al. (6.75 Gbp  
181 to 7.63 Gbp, (2011)), Ochatt et al. (7.82 to 8.90 Gbp, (2013)), Nandini et al. (6.85 Gbp, (1997)) and  
182 Macas et al. (6.52 Gbp (2015)). Legume genomes, particularly species in the Viciaea (=Fabeae)  
183 including the genus *Lathyrus*, are highly variable in size, and this variation correlates strongly with  
184 the copy numbers of repeated elements (Leitch et al. 2019; Vondrak et al. 2020). The reported  
185 variation in genome size reported for *Lathyrus sativus* could be due to intra-specific variation  
186 between genotypes or experimental error. The range of variation in genome size for *Pisum sativum*  
187 in the Kew database is 3.49 to 5.42, which is a greater range than for grass pea, and an even greater  
188 proportional range (Leitch et al. 2019). As our genome size estimate for LS007 depends on the size of  
189 the *Pisum sativum* genome which we used as a standard, the relation between our estimate and the  
190 sizes of the assemblies should not be over-interpreted.

191 OGRE elements are a type of Ty3/gypsy LTR retrotransposon first discovered in legumes (Neumann,  
192 Požárková, and Macas 2003; Neumann et al. 2006), but since found in other plant families as well  
193 (Macas and Neumann 2007). OGRE elements are characterised by their large size (up to 25 kbp) and  
194 the presence of an additional ORF upstream of the gag-pro-pol polyprotein ORF usually present in  
195 LTR retrotransposons (Macas and Neumann 2007). In a survey of the genomes of 23 species within  
196 the legume family Viciaea (=Fabeae), OGRE elements typically make up about 40% of the entire  
197 genome (22.5 - 64.7%), and OGRE-content correlates strongly with genome size (Macas et al. 2015).  
198 Our sequencing data gave an estimated OGRE-content of 37% for the LS007 nuclear genome. The  
199 repeat analysis using RepeatExplorer2 (Macas et al. 2015) used by both ourselves and Macas et al.  
200 relies on <1X coverage raw short-read data and is therefore unaffected by the difficulties of  
201 assembling repeat regions that often lead to repeat regions being under-represented in genome  
202 assemblies. Nevertheless, our result of 37.3% OGRE-content in LS007 contrasts with the 45.5%  
203 estimate reported for *Lathyrus sativus* by Macas et al. (2015). This may be due to genotypic  
204 differences between LS007 and the commercial line used by Macas et al.

205 Satellites in *Lathyrus sativus* have been claimed to originate from the LTRs of OGRE elements  
206 (Vondrak et al. 2020). Rapid tandem duplication of such repeated elements may be a mechanism for  
207 the high degree of species specificity of satellite elements in legumes. A single satellite element  
208 (FabTR53\_LAS\_A) with a consensus sequence of 660 bp and no significant blastn hits against the  
209 NCBI\_nt database represented 4.7% of our short reads, corresponding to an estimated 289.8 Mbp or  
210 440,000 copies across the genome. One assembly scaffold carried 144 copies of this satellite in  
211 tandem, spread over a 88kbp region. Fluorescence in-situ hybridisation (FISH) conducted by Vondrak  
212 et al. (2020) showed that this satellite is concentrated in the sub-telomeric regions of the  
213 chromosomes, rather than the primary constrictions of chromosomes, as is typical for most  
214 satellites. This high level of repetition underlines the difficulty of assembling the grass pea genome  
215 and the necessity for long sequence reads that can span repetitive regions.

216 Consequently, the assembly of the grass pea genome from short paired-end read data alone reached  
217 only a contig N50 of 5.7 kbp, as many repeated regions could not be adequately resolved.  
218 Scaffolding with Illumina long mate pair data derived from 2 kbp, 5 kbp, 8 kbp and 14 kbp libraries  
219 was able to link contigs into scaffolds reaching an N50 of 59.7 kbp. In the process of scaffolding,  
220 many sections of unknown sequence originating from LMP inserts were introduced into the  
221 scaffolds, resulting in a total of 1.9 billion Ns in the Eiv1 assembly. It is possible that some smaller  
222 contigs may be nested within the unknown sequence of these scaffolds, resulting in a true assembly  
223 shorter than the total of 8.12 Gbp. Technical limitations prevent the generation of substantially  
224 longer mate pair inserts with reliable sizes. Improving this assembly thus required the use of  
225 alternative sequencing and/or scaffolding technologies.

226 To overcome the challenges of size and repetitiveness of the grass pea genome, we opted for a  
227 hybrid assembly approach, combining the high yields and base accuracy of Illumina sequencing with  
228 the long reads of Oxford Nanopore sequencing technologies, an approach recommended by Kreplak  
229 et al. (2019). Long-read sequencing allowed us to increase assembly contiguity compared to the  
230 long-mate-pair-scaffolded Illumina paired-end assembly (155 kbp contig N50 vs. 59 kbp scaffold  
231 N50). In addition to the Redbean assembly we report here, we ran Flye (Kolmogorov et al. 2019),  
232 which produced a 14 Gbp assembly (before crashing during polishing), and Shasta (Shafin et al.  
233 2019), which produced only a 169 Mbp assembly (results not shown), using the same PromethION  
234 data. Neither of these were pursued further. We also ran MaSuRCA (Zimin et al. 2013) using both  
235 PromethION and Illumina paired-end data, but the algorithm was unable to complete due to the size  
236 of the input datasets.

237 BUSCO completeness scores were consistently higher in the polished Redbean assembly across 4  
238 BUSCO databases (Simão et al. 2015), although none reached >90% completeness. The polished  
239 Redbean assembly exhibited significantly lower rates of duplicated BUSCOs than the Elv1 assembly,  
240 on average 4.2% vs. 14.6%. This might reflect duplication of pseudo-hemizygous regions caused by  
241 heterozygosity in the sequenced genotype. Although the genotype we used underwent single-seed  
242 descent, it is possible that some heterozygotic regions remained in the sequenced material. The  
243 longer PromethION reads used by the Redbean assembler collapsed more of these heterozygous  
244 regions into single contigs, suggesting that these regions were indeed heterozygous. Duplicated  
245 pseudo-hemizygous regions can also be further collapsed using the Purge\_dups tool (Guan et al.  
246 2020), but we elected not to do this due to the risk of collapsing truly duplicated genomic regions.

247 The annotation of the polished Redbean assembly using the transcriptome datasets already  
248 described is pending and this draft will be amended accordingly, as soon as it is ready.

249 The assembly and annotation datasets are available to researchers on request on receipt of a  
250 written commitment not to publish analyses of the data before publication of the Consortium paper  
251 or on the basis of agreement with all the authors.

252

## 253 **Materials and Methods**

### 254 **Genome size estimation by flow cytometry**

255 Grass pea genome size was estimated following the procedure described by Dolezel et al. (2007).  
256 Fresh, young leaf tissue (40 mg) of grass pea (LS007) and *Pisum sativum* (semi-leafless, obtained  
257 from a local market in Nairobi) was sliced finely using a scalpel blade while immersed in 2 mL of ice-  
258 cold Galbraith buffer (45 mM MgCl<sub>2</sub>, 30 mM sodium citrate, 20 mM MOPS, 0.1% w/v Triton X-100,  
259 pH 7). Three biological replicates were prepared for each grass pea and pea. Supernatants were  
260 filtered through one layer of Miracloth (pore size 22-25 µm). One aliquot of 600 µL was prepared  
261 from each replicate, along with three grass pea + pea mixes at 2:1, 1:1 and 1:2 ratios, respectively.  
262 Propidium iodide was added to each tube to a concentration of 50 µM. Reactions were incubated for  
263 1 h on ice before measuring on a FACSCantoll flow cytometer (Becton Dickinson) with flow rates  
264 adjusted to 20-50 events/s. Results were analysed using FCSalyser (v. 0.9.18 alpha). Grass pea  
265 genome size was estimated from the mixed sample by dividing the mean of the PE-A peak  
266 corresponding to grass pea nuclei by the mean of the PE-A peak corresponding to pea nuclei and  
267 multiplying by 4.300 Gbp, the estimated genome size of pea (Leitch et al. 2019).

### 268 **Illumina Sequencing**

#### 269 **Genomic DNA isolation**

270 Seeds of grass pea (*Lathyrus sativus*) were obtained from King's of Coggershall and underwent six  
271 generations of single-seed-descent at the John Innes Centre Germplasm Research Unit. This  
272 accession, named LS007 is of European origin, white-flowered, with fully cream-coloured, large and  
273 flattened seeds. LS007 seeds are available from the JIC GRU. Seeds were surface sterilised and  
274 germinated on distilled water agar in the dark. Genomic DNA was isolated from the etiolated  
275 seedlings using a modified CTAB protocol and subsequently, high molecular weight DNA purified  
276 using the Qiagen MagAttract kit.

#### 277 **Initial Paired end (PE) library DNA sequencing and assembly**

278 Paired end (PE) sequencing was carried out using PCR-free libraries on the HiSeq Illumina platform.  
279 Five lanes of untrimmed PE data from PSEQ-907 (LIB21060) were assembled using both Abyss  
280 (Simpson et al. 2009) and Discover de novo v.52488 (Love et al. 2016). A range of kmer spectra plots  
281 were generated using KAT (Mapleson et al. 2017) to determine how the kmers from the initial PE  
282 library were represented in these final PE assemblies. Given broad similarity between kmer spectra



283 for both assemblies, and that the best assembly statistics were returned by Discover de novo (best  
284 Abyss N50 = 1.5 kbp, best Discover de novo N50 = 5.7 kbp), the Discover de novo assembly was  
285 selected as the best candidate for further improvement/scaffolding. The total usable coverage from  
286 clipped LMP libraries (using categories A, B and C as identified by *Nextclip*) are given in Supplemental  
287 table S3

## 288 Sequencing and assembly of Long Mate Pair (LMP) libraries

289 A pool of Long Mate Pair (LMP) libraries with a range of insert sizes (12 in total) was prepared and  
290 sequenced. These reads were processed with *nextclip*, and mapped to both the best Abyss and  
291 Discover de novo assemblies. Binary Alignment Map (BAM) files were sorted and de-duplicated  
292 before running the Picard tool *CollectInsertSizeMetrics* on them (Wysokar et al. 2016). The Discover  
293 de novo contigs/scaffolds were generally longer and therefore judged to be the better set for  
294 mapping, particularly for the LMPs with longer insert sizes. The insert size estimates from  
295 *CollectInsertSizeMetrics* were supplemented with additional statistics from the Picard  
296 *JumpingLibraryMetrics* module. Analysis of both these sets of metrics suggested that the LIB21840  
297 (~2 kbp), LIB21836 (~5 kbp) and LIB21834 (~8 kbp) LMPs were the three best libraries for  
298 resequencing for greater coverage depth.

299 These three LMP libraries were pooled (IPO3823) and sequenced in four separate lanes across three  
300 different sequencing runs (PSEQ-1057, PSEQ-1067 and PSEQ-1101). *Nextclip* was used on the reads  
301 after each run to assess the extent to which usable LMP coverage had been recovered.

## 302 Scaffolding

303 After an initial scaffolding attempt using SSPACE v2.0 (Boetzer et al. 2011), an extra lane of the  
304 longest 8 kbp LMP (LIB21834) was sequenced on PSEQ-1159, and this helped to increase the final  
305 clipped coverage for this library to 10.7x (assuming a genome size of 6.3 Gbp). Use of the Bowtie  
306 alignment option of SSPACE (v3.0) was rapidly judged to be the only feasible option (because Bwa  
307 was extremely slow), but this necessitated making a custom installation of SSPACE which was able to  
308 handle genome references of >4 Gbp. This version of SSPACE was run specifying k=3 as the  
309 minimum number of connections required before scaffolding two sequences together. All  
310 contigs/scaffolds of >500 bp from the Discover de novo assembly were provided as the input  
311 reference. The final scaffold N50 using 500 bp as the input length cutoff was ~32 kbp. This assembly  
312 was re-scaffolded using an even longer insert LMP (LIB28873, ~14 kbp insert size). Using exactly the  
313 same scaffolding procedure and parameters, providing only the clipped LIB28873 reads (final  
314 coverage is 5.2x assuming 6.3 Gbp genome size), the final scaffold N50 was ~47 kbp. Increasing the

315 minimum length to 1 kbp improved the assembly, resulting in a final draft assembly of 8.12 Gbp  
316 (comprising 6.2 Gbp known nucleotides and 1.9 Gbp Ns) and a scaffold N50 of 59.7 kbp. The final  
317 assembly statistics are shown in Table 2 and Figure 4.

318 Annotation of the Illumina assembly (Elv1)The gene annotation pipeline used to annotate the LS007  
319 Elv1 genome assembly has been adapted from the pipeline used to annotate the wheat genome  
320 developed at the Earlham Institute (Venturini, Kaithakottil, and Swarbreck 2016).

### 321 Identification of repeats from the genome assembly

322 RepeatModeler (v1.0.10 - <http://www.repeatmasker.org/RepeatModeler/>) was used for *de novo*  
323 identification of repetitive elements from the assembled grass pea genome sequence. Protein  
324 coding genes in the RepeatModeler generated library were hard-masked (i.e. replaced with Ns) using  
325 the Arabidopsis Araport11 dataset and *Cicer arietinum* Annotation 101  
326 ([https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/Cicer\\_arietinum/101/](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Cicer_arietinum/101/)) coding genes. Any  
327 genes with descriptions indicating “transposon” or “helicase” were removed. TransposonPSI  
328 (r08222010) (Haas 2010) was run and significant hits hard-masked, with the output being used to  
329 mask the RepeatModeler library. Unclassified repeats were searched in a custom BLAST database of  
330 organellar genomes (mitochondrial and chloroplast sequences from Fabales [ORGN] in NCBI  
331 nucleotide division, downloaded on 22.9.2017). Any repeat families matching organellar DNA were  
332 also hard-masked.

333 Repeat identification was refined by running RepeatMasker v4.0.7 (Smit, Hubley, and Green 2013)  
334 with RepBase Viridiplantae library and with the customized repeatmodeler library (i.e. after masking  
335 out protein coding genes), both using the -nolow setting.

336 The combined masking resulted in ~62% of the grass pea genome being soft-masked (i.e. rendered  
337 in lowercase to stop alignment of transcriptome data).

338

### 339 Reference guided transcriptome reconstruction

#### 340 Alignment of Illumina RNA-Seq data

341 RNA-Seq data from three different genotypes LS007, LSWT11 and Mahateora (Emmrich 2017;  
342 Emmrich et al. 2019) was used for grass pea genome annotation: 12 samples from LS007 and  
343 Mahateora each (each root and shoot tissues from 3 biological replicates of droughted and non-  
344 droughted treatments) and 7 samples from seedling shoot tip, seedling root tip, whole root, whole

345 leaves, early flowers, early pods and late pods from LSWT11 (Table 2) for a total of 31 individual  
346 RNA-seq samples). In total, the three filtered datasets comprised over 2.5 billion paired-end reads.  
347 For each dataset, read samples were collapsed by tissue and filtered using Xtool (Nazar et al. 2020),  
348 with the command line options. Due to concerns about high concentrations of ribosomal RNA,  
349 datasets were further filtered using SortMeRNA v. 2.0 (Kopylova, Noé, and Touzet 2012), and using  
350 RFam (5S and 5.8S) and Silva (Archaea 16S-23S, Bacteria 16S-23S, Eukaryota 18S-28S) as databases.

### 351 Alignment with HISAT2

352 Filtered reads were aligned to the grass pea genome using HISAT2 v2.0.5 (Kim, Langmead, and  
353 Salzberg 2015, 2) with options: --min-intronlen=20 --max-intronlen=50000 --dta-cufflinks. The RNA-  
354 seq data and alignments are summarised in supplemental table S4.

### 355 Transcript assembly

356 The Illumina RNA-Seq alignments (31 RNA-Seq transcript alignments from HISAT2) were assembled  
357 using two different RNA-Seq assembly tools: Cufflinks v. 2.2.1 (Trapnell et al. 2010; Roberts et al.  
358 2011), with options: --min-intron-length=20 --max-intron-length=50000 and StringTie v.1.3.3 (Pertea  
359 et al. 2015), with options: -f 0.05 -m 200. The results of transcript assembly are shown in  
360 supplemental table S5.

361 Mikado v1.0.1 (Venturini et al. 2018), which uses transcript assemblies generated by multiple  
362 methods to improve transcript reconstruction, was used to integrate the ~4.5 million Illumina  
363 contigs generated using Cufflinks and StringTie (see supplemental table S6). Loci were first defined  
364 across all assemblies, followed by scoring of each transcript for ORF and cDNA length, position of  
365 the ORF within the transcript (and presence of multiple ORFs), UTR lengths. The highest scoring  
366 transcript assembly was used along with any additional transcripts (splice variants) compatible with  
367 this representative transcript assembly. Mikado selected transcript sets were generated for use in  
368 gene predictor training and annotation (Table 4) using the RNA-Seq transcript assemblies with the  
369 “chimera\_split” option set to PERMISSIVE. For Mikado runs incorporating BLAST data, transcripts  
370 that had passed the “prepare” step were compared to filtered and masked proteins from *Cicer*  
371 *arietinum*, *Cucumis sativus*, *Fragaria vesca*, *Glycine max*, *Malus domestica*, *Medicago truncatula*,  
372 *Prunus persica*, *Phaseolus vulgaris*, *Trifolium pratense* using BLAST+ v. 2.6.0 (Altschul et al. 1997).  
373 Each result was limited to the best 15 matches.

## 374 Gene prediction using evidence guided AUGUSTUS

375 Protein coding genes were predicted using AUGUSTUS (Stanke et al. 2004; 2006), which uses a  
376 Generalized Hidden Markov Model (GHMM) employing both intrinsic and extrinsic information.

### 377 *Intron/exon junctions*

378 RNA-Seq junctions were derived from RNA-Seq alignments using Portcullis v1.0.2 (Mapleson et al.  
379 2018) with default filtering parameters. Junctions that passed the Portcullis filter with a score of 1  
380 were classified as 'Gold' and with a score <1 were classified as 'Silver'.

### 381 *Proteins*

382 Predicted protein sequences from 9 species (*Cicer arietinum*, *Cucumis sativus*, *Fragaria vesca*,  
383 *Glycine max*, *Malus domestica*, *Medicago truncatula*, *Prunus persica*, *Phaseolus vulgaris*, *Trifolium*  
384 *pratense*) were soft-masked for low complexity regions (segmasker from NCBI BLAST+ 2.6.0) and  
385 aligned to the soft-masked grass pea genome sequence (using repeatmodeler repeats) with  
386 exonerate v2.2.0 (Slater and Birney 2005) with parameters: --model protein2genome --  
387 showtargetgff yes --showvulgar yes --softmaskquery yes --softmasktarget yes --bestn 10 --minintron  
388 20 --maxintron 50000 --score 50 --geneseed 50.

389 To identify high confidence alignments, exonerate results were filtered at 50% identity and 80%  
390 coverage. Alignments with introns longer than 10 kbp were removed from further analyses (see  
391 Supplemental table S7).

### 392 *Classification of Mikado transcripts:*

393 The primary Mikado models for each locus were classified into three categories:

Gold	Mikado transcripts having a full length hit (complete/putative complete) with Full-LengtherNEXT v0.0.8 (Seoane, Fernandez, and Guerrero 2018) and a maximum of 2 complete 5'UTR exons and 1 complete 3'UTR exon consistent Full-lengtherNEXT and TransDecoder v2.0.1 (Grabherr et al. 2011) CDS coordinates, a minimum CDS to transcript ratio of 50% and a single transcript per gene.
Silver	Remaining models meeting UTR restrictions with the additional constraint of having a CDS length of at least 900bp
Bronze	Any remaining Mikado transcripts were assigned to bronze

## 394 Gene prediction

### 395 Gene predictor training

396 The Mikado gold set transcripts were selected for training AUGUSTUS (Grabherr et al. 2011). We  
397 excluded genes with a genomic overlap within 1000bp of a second gene and gene models that were  
398 homologous to each other with coverage and identity  $\geq 80\%$ . The filtered set contained 9604  
399 transcripts, from which 2000 transcripts were selected at random for training AUGUSTUS and  
400 another 200 transcripts were used for testing. The trained AUGUSTUS model resulted in values of  
401 0.976 sensitivity (sn), 0.909 specificity (sp) at the nucleotide level, sn 0.837, sp 0.801 at the exon  
402 level and sn 0.41, sp 0.376 at the gene level.

403 AUGUSTUS (v2.7) was used to predict gene models for the Elv1 assembly using the evidence hints  
404 generated from nine sets of cross species protein alignments (listed above), Mikado Illumina models  
405 and intron/exon junctions defined using the RNA-Seq data . Interspersed repeats were provided as  
406 “nonexonpart” to exclude them from analysis. We assigned additional bonus scores and priority  
407 based on evidence type and classification (Gold, Silver, Bronze) to reflect the reliability of different  
408 evidence sets.

### 409 Cross species protein similarity ranking

410 Each gene model was assigned a protein rank (P1–P5) reflecting the level of coverage of the best  
411 identified homolog in plant protein databases. Protein ranks were assigned according to:

---

Protein Rank 1 (P1)	proteins identified as full length in Full-LengtherNEXT in the UniProt database or at least 90% coverage in a supplementary BLAST database consisting of proteins from: <i>Cicer arietinum</i> , <i>Cucumis sativus</i> , <i>Fragaria vesca</i> , <i>Glycine max</i> , <i>Malus domestica</i> , <i>Medicago truncatula</i> , <i>Prunus persica</i> , <i>Phaseolus vulgaris</i> , and <i>Trifolium pratense</i> proteins
Protein Rank 2 (P2)	proteins with at least 60–90% coverage in the supplementary BLAST database
Protein Rank 3 (P3)	proteins with at least 30–60% coverage in the supplementary BLAST database
Protein Rank 4 (P4)	proteins with a low coverage hit (between 0–30%) in the supplementary BLAST database
Protein Rank 5 (P5)	proteins with no hit in the supplementary BLAST database

---

## 412 Grass pea transcript support ranking

413 Transcript ranks (T1–T5) were assigned based on the support for each predicted gene model from  
414 the assembled grass pea RNA-Seq data (all 4,582,819 transcripts assembled using Cufflinks and  
415 StringTie).

416 We calculated a variant of annotation edit distance (*AED*) and used this to determine a transcript  
417 level ranking.

$$418 \text{AED} = 1 - (SN + SP)/2$$

419 where *SN* is sensitivity and *SP* specificity.

420 *AED* was calculated at base, exon and splice junction levels against all individual transcripts used in  
421 our gene build (Illumina assemblies). The mean of base, exon and junction *AEDs* (derived using the  
422 Mikado ‘compare utility’) based on the transcript that best supported the gene model was then used  
423 to assign transcript rank as shown:

<b>Transcript Rank 1 (T1)</b>	Full length support single cDNA reads (not relevant to this study, but retained for workflow compatibility)
<b>Transcript Rank 2 (T2)</b>	full length support from Illumina assemblies
<b>Transcript Rank 3 (T3)</b>	Best average <i>AED</i> less than 0.5
<b>Transcript Rank 4 (T4)</b>	Best average <i>AED</i> between 0.5 and 1
<b>Transcript Rank 5 (T5)</b>	No transcriptomic support (best average <i>AED</i> = 1)

## 424 Assignment of gene biotypes and confidence classification

425 Gene models were classified as either coding or repeat associated and classified as high- or low-  
426 confidence based on cross-species protein similarity and grass pea transcript evidence. Each  
427 transcript was assigned a ‘protein rank’ and a ‘transcript rank’ (as described above) and a binary  
428 confidence tag (“High” vs “Low”), combining the two. Any transcript with protein rank P2 or better  
429 as well as transcript rank T4 or better were classified as ‘high confidence’, all others as ‘low  
430 confidence’.

431

## 432 Assignment of a locus biotype

433 Using the protein and transcript rankings, we assigned a locus biotype (‘repeat-associated’ or  
434 ‘protein-coding’) to each gene.

435 Repeat-associated biotypes

436 Genes were classified as 'repeat-associated' if ANY of their transcripts: :

- 437 - aligned with at least 20% similarity and 30% coverage to the TransposonPSI library
- 438 v08222010 (Haas 2010)
- 439 - or had at least 30% coverage by RepeatModeler/RepeatMasker derived interspersed repeats
- 440 - or had a match for "retrotransposon", "transposon", "helicase" in their annotation (AHRD
- 441 and interproscan) description.

442 *Protein-coding genes*

443 Genes not assigned as 'repeat-associated', were assigned the 'protein coding' biotype, along with all  
444 transcripts associated with them.

445 **Removal of spurious genes**

446 After assigning biotypes to each gene, we used Kallisto v 0.43.1 (Bray et al. 2016) to estimate gene  
447 expression in all of our samples. Genes were retained if one of their transcripts had either an  
448 expression level exceeding 0.5 transcripts per million (TPM) in at least one sample, as measured by  
449 Kallisto OR any BLAST hits from the Full-LengtherNEXT analysis against the UniProt database.

450 Any gene whose transcripts were all marked for removal was excluded from the final annotation.

451 **Assignment of a representative gene model**

452 We assigned representative gene models by selecting the model with the highest confidence ranking  
453 (as shown in Supplemental table S2), and the lowest *AED* with the order of priority for ranking:

- 454 1. highest protein rank
- 455 2. highest transcript rank
- 456 3. lowest *AED*.

457 **Functional annotation of protein coding transcripts**

458 Proteins were annotated using AHRD v.3.1 (Hallab 2014). Predicted protein sequences were  
459 searched against Araport11 *A. thaliana* protein sequences (Cheng et al. 2017) and the plant  
460 sequences of UniProt v. 23Nov2017, both SwissProt and TrEMBL datasets (The UniProt Consortium  
461 2019) using BLASTP+ v. 2.6.0 asking for a maximum e-value of 1e-5. We also ran InterProScan  
462 5.22.61 (Jones et al. 2014) and provided the InterProScan output to AHRD. We adapted the standard  
463 example configuration file pathstest/resources/ahrd\_example\_input.yml, distributed with the AHRD

464 tool by providing the GOA mapping from UniProt. Supplemental table S2 reports the final number of  
465 coding transcripts per each rank.

466

## 467 PromethION Sequencing

468 Two methods of DNA extraction and three methods of DNA processing were used to optimise yield  
469 and read length profile of the PromethION sequencing.

### 470 DNA extraction

#### 471 *Circulomics Nanobind Plant Nuclei Kit extraction (CN)*

472 Fresh, aseptic LS007 shoot tissue (0.5 g) was ground in liquid nitrogen using a mortar and pestle. The  
473 resulting powder was resuspended by vortexing in 10 mL of nuclei extraction buffer (0.35 M Sorbitol,  
474 100 mM Tris pH 8.0, 5mM EDTA and 1%  $\beta$ -mercaptoethanol) and then filtered through Miracloth.  
475 The filtrate was centrifuged at 750 x g for 15 min at 4 °C and the pellet was resuspended in nuclei  
476 extraction buffer (with the addition of 0.4% Triton). The wash was repeated before centrifugation  
477 and resuspension in 1mL of nuclei resuspension buffer (0.35 M Sorbitol, 100 mM Tris, 5mM EDTA).  
478 The suspension was then centrifuged at 750 x g for 15 min at 4 °C and the supernatant was removed.  
479 The resulting nuclear pellet was taken into the Circulomics Nanobind Plant nuclei Kit (Circulomics;  
480 SKU NB-900-801-01) according to the manufacturers protocol. The resulting DNA had a peak  
481 fragment size >60 kbp.

482

#### 483 *Qiagen DNeasy Plant Mini Kit extraction (QD)*

484 0.5 g of fresh, aseptic LS007 shoot tissue was ground under liquid nitrogen using a mortar and pestle  
485 and resuspended in 2 mL of AP1 (Qiagen) with 20  $\mu$ l of RNase I (Qiagen). This was incubated at 65°C  
486 for 10 minutes, before aliquoting into 5 tubes and proceeding with the Qiagen DNeasy Plant Mini Kit  
487 (Qiagen; 69104) protocol. Eluted DNA was pooled into a single tube and had a peak fragment size of  
488 48 kbp.

### 489 DNA processing

490 The process of Sample preparation was iteratively optimized to achieve both high total sequence  
491 yield and long read length. The exact combination of techniques used for each library preparation is  
492 given in Table 1.



493 *Ampure XP bead purification/concentration (AMP)*

494 To increase the concentration of DNA for input into the Short Read Eliminator and to reduce the loss  
495 of DNA in library preparation, All except the DNA Qiagen DNeasy-extracted samples were pre-  
496 incubated with 0.8x Ampure XP (Beckman Coulter; A63881). Beads were washed twice in 80%  
497 ethanol and DNA was eluted from beads with TE buffer (10 mM Tris HCl pH 8.0, 1 mM EDTA).

498 *Needle shearing (NS)*

499 Starting with the 2<sup>nd</sup> loading of flowcell FC1, all except the Qiagen DNeasy-extracted samples were  
500 needle-sheared 30 times with a 26 gauge needle to reduce the amount of very high molecular  
501 weight DNA, which can cause blocking and therefore an artificially low N50.

502 *Short Read Elimination (SRE)*

503 Starting with the 3<sup>rd</sup> loading of flowcell FC1, all samples were subjected to a Circulomics Short Read  
504 Eliminator (Circulomics; SKU SS-100-101-01) treatment to reduce the number of short fragments.

505 **Library preparation and loading**

506 All libraries were prepared using the Genomic DNA by Ligation (SQK-LSK109) – PromethION Kit  
507 (Oxford Nanopore Technologies; SQK-LSK109) following the manufacturer’s procedure. Libraries  
508 were loaded onto PromethION Flow Cells (Oxford Nanopore Technologies; FLO-PRO002) at between  
509 250-400 ng. Due to the rapid accumulation of blocked flow cell pores or due to apparent read length  
510 anomalies on some runs, flow cells used in runs were treated with a nuclease flush to digest blocking  
511 DNA fragments before reloading with fresh library according the Oxford Nanopore Technologies  
512 Nuclease Flush protocol, version NFL\_9076\_v109\_revD\_08Oct2018.

513

514 **Data processing**

515 Fast5 sequences produced by PromethION sequencing were basecalled using the Guppy high  
516 accuracy basecalling model (dna\_r9.4.1\_450bps\_hac.cfg) and the resulting fastq files were quality  
517 filtered by the basecaller. Fastq files from all five sequencing runs were pooled and assembled using  
518 Redbean (Ruan and Li 2020, 2) (previously wtdbg2, version 2.2). This assembly was polished with  
519 minimap2 (v. 2.17) (Li 2018, 2) using the original nanopore dataset in fasta format, followed by  
520 polishing using bwa (v. 0.7.17) (Li 2013).

## 521 BUSCO analysis

522 Gene space completeness in Rbp and Elv1 was assessed using BUSCO v.4.0.4\_cv1 (Simão et al. 2015)  
523 and the odb10 databases for eukaryota, viridiplantae, eudicots and fabales, employing default  
524 parameters.

## 525 Phylogeny

526 Ten arbitrary single-copy BUSCO sequences found in the Rbp and Elv1 assemblies and their  
527 predicted mRNA homologs with the highest NCBI BLAST scores in *Arabidopsis thaliana*, *Cajanus*  
528 *cajan*, *Cicer arietinum*, *Glycine max*, *Lotus japonicus*, *Lupinus angustifolius*, *Medicago truncatula*,  
529 *Nelumbo nucifera*, *Phaseolus vulgaris*, *Pisum sativum*, *Populus trichocarpa*, *Prunus persica*, *Solanum*  
530 *lycopersicum*, *Theobroma cacao*, *Vigna angularis*, *Vigna radiata* and *Vitis vinifera* (Sayers et al.  
531 2019). Multiple sequence alignments for each homolog group were produced using MUSCLE (Edgar  
532 2004). We merged the 10 unrooted trees each containing 19 unique taxa (including the two LS007  
533 assemblies, Elv1 and Rbp) with 79 unique splits from total of 360, using DendroPy 3.12.0 (Sukumaran  
534 and Holder 2010), and the consensus tree was built using the minimum clade frequency threshold  
535 set at 0.5.

## 536 Repeat analysis

537 Genome repeat structure was analysed using RepeatExplorer2 (Macas et al. 2015), a graph-based  
538 read clustering algorithm for repeat identification. Illumina HiSeq paired end data was downsampled  
539 to a coverage of 0.1X and trimmed, quality filtered, cutadapt filtered and interlaced into a single  
540 FASTA file for processing by repeatExplorer2 using the ELIXIR CZ Galaxy server using default  
541 parameters (Afgan et al. 2016).

542

543

## 544 Acknowledgments

545 This work was supported by a John Innes Centre Institute Development Grant, the Biotechnology  
546 and Biological Sciences Research Council (BBSRC) Detox Grass pea project (BB/L011719/1) the BBSRC  
547 SASSA UPGRADE project (BB/R020604/1), the BBSRC Institute Strategic Programme  
548 (BBS/E/J/000PR9799) and the Nottingham Future Food Beacon of Excellence. PMFE's studentship  
549 was funded by the John Innes Foundation's Student Rotation programme. None of the funding  
550 bodies were involved in the design of this study, the collection or analysis or interpretation of data,  
551 or in writing the manuscript.

552 The Galaxy server that was used for some calculations is in part funded by Collaborative Research  
553 Centre 992 Medical Epigenetics (DFG grant SFB 992/1 2012) and German Federal Ministry of  
554 Education and Research (BMBF grants 031 A538A/A538C RBC, 031L0101B/031L0101C de.NBI-epi,  
555 031L0106 de.STAIR (de.NBI)).

556 We acknowledge Guru Radhakrishnan, Tjelvar Olsen, Matt Hartley, Shabhonam Caim, Pirita  
557 Paaajanen and Burkhard Steuernagel for their vital support and advice in bioinformatics and data  
558 handling.

559

## 560 Availability of data and materials

561 All plant materials used in this article are available from the corresponding author. The datasets  
562 used and analysed during the current study have been uploaded (under embargo) to the European  
563 Nucleotide Archive but are available from the corresponding author upon reasonable request. All  
564 data will become publicly available upon publication of the peer-reviewed article.

565

566

## 567 References

- 568 Afgan, Enis, Dannon Baker, Marius van den Beek, Daniel Blankenberg, Dave Bouvier, Martin Čech,  
569 John Chilton, et al. 2016. 'The Galaxy Platform for Accessible, Reproducible and  
570 Collaborative Biomedical Analyses: 2016 Update'. *Nucleic Acids Research* 44 (Web Server  
571 issue): W3–10. <https://doi.org/10.1093/nar/gkw343>.
- 572 Altschul, Stephen F, Thomas L Madden, Alejandro A Schäffer, Jinghui Zhang, Zheng Zhang, Webb  
573 Miller, and David J Lipman. 1997. 'Gapped BLAST and PSI-BLAST: A New Generation of  
574 Protein Database Search Programs'. *Nucleic Acids Research* 25 (17): 3389–3402.
- 575 Boetzer, Marten, Christiaan V. Henkel, Hans J. Jansen, Derek Butler, and Walter Pirovano. 2011.  
576 'Scaffolding Pre-Assembled Contigs Using SSPACE'. *Bioinformatics* 27 (4): 578–79.  
577 <https://doi.org/10.1093/bioinformatics/btq683>.
- 578 Bray, Nicolas L., Harold Pimentel, Páll Melsted, and Lior Pachter. 2016. 'Near-Optimal Probabilistic  
579 RNA-Seq Quantification'. *Nature Biotechnology* 34 (5): 525–27.  
580 <https://doi.org/10.1038/nbt.3519>.
- 581 Campbell, Clayton Garnet. 1997. *Grass Pea: Lathyrus Sativus L. Promoting the Conservation and Use*  
582 *of Underutilized and Neglected Crops*. Vol. 18. International Plant Genetic Resources  
583 Institute.
- 584 Chakrabarti, A, I M Santha, and S L Mehta. 1999. 'Molecular Characterisation of Low ODAP  
585 Somaclones of Lathyrus Sativus'. *Journal of Plant Biochemistry and Biotechnology* 8 (1): 25–  
586 29.
- 587 Challis, Richard. (2015) 2020. *Rjchallis/Assembly-Stats*. JavaScript.  
588 <https://github.com/rjchallis/assembly-stats>.
- 589 Cheng, Chia-Yi, Vivek Krishnakumar, Agnes P. Chan, Françoise Thibaud-Nissen, Seth Schobel, and  
590 Christopher D. Town. 2017. 'Araport11: A Complete Reannotation of the Arabidopsis  
591 Thaliana Reference Genome'. *The Plant Journal* 89 (4): 789–804.  
592 <https://doi.org/10.1111/tpj.13415>.
- 593 Cohn, D F, and M Streifler. 1983. 'Intoxication by the Chickling Pea (Lathyrus Sativus): Nervous  
594 System and Skeletal Findings'. In *Toxicology in the Use, Misuse, and Abuse of Food, Drugs,*  
595 *and Chemicals*, 190–93. Springer.
- 596 Dolezel, Jaroslav, Johann Greilhuber, and Jan Suda. 2007. 'Estimation of Nuclear DNA Content in  
597 Plants Using Flow Cytometry'. *Nature Protocols* 2 (9): 2233–44.  
598 <https://doi.org/10.1038/nprot.2007.310>.
- 599 Drouin, Pascal, Danielle Prévost, and Hani Antoun. 2000. 'Physiological Adaptation to Low  
600 Temperatures of Strains of Rhizobium Leguminosarum Bv. Viciae Associated with Lathyrus  
601 Spp'. *FEMS Microbiology Ecology* 32 (2): 111–20.
- 602 Dufour, D L. 2011. 'Assessing Diet in Populations at Risk for Konzo and Neurolathyrism'. *Food and*  
603 *Chemical Toxicology* 49 (3): 655–61. <https://doi.org/10.1016/j.fct.2010.08.006>.
- 604 Edgar, Robert C. 2004. 'MUSCLE: Multiple Sequence Alignment with High Accuracy and High  
605 Throughput'. *Nucleic Acids Research* 32 (5): 1792–97. <https://doi.org/10.1093/nar/gkh340>.
- 606 Emmrich, Peter M. F. 2017. 'Genetic Improvement of Grass Pea (Lathyrus Sativus) for Low  $\beta$ -L-ODAP  
607 Content'. University of East Anglia. <https://ueaeprints.uea.ac.uk/63944/>.
- 608 Emmrich, Peter M. F., Martin Rejzek, Lionel Hill, Paul Brett, Anne Edwards, Abhimanyu Sarkar, Rob A.  
609 Field, Cathie Martin, and Trevor L. Wang. 2019. 'Linking a Rapid Throughput Plate-Assay with  
610 High-Sensitivity Stable-Isotope Label LCMS Quantification Permits the Identification and  
611 Characterisation of Low  $\beta$ -L-ODAP Grass Pea Lines'. *BMC Plant Biology* 19 (1): 489.  
612 <https://doi.org/10.1186/s12870-019-2091-5>.
- 613 Ghasem, Karimzadeh, Maryam Danesh-Gilevaei, and Majid Aghaalikhani. 2011. 'Karyotypic and  
614 Nuclear DNA Variations in Lathyrus Sativus (Fabaceae)'. *Caryologia* 64 (1): 42–54.  
615 <https://doi.org/10.1080/00087114.2011.10589763>.

- 616 Girma, A, B Tefera, and L Dadi. 2011. 'Grass Pea and Neurolathyrism: Farmers' Perception on Its  
617 Consumption and Protective Measure in North Shewa, Ethiopia'. *Food and Chemical*  
618 *Toxicology* 49 (3): 668–72. <https://doi.org/10.1016/j.fct.2010.08.040>.
- 619 Grabherr, Manfred G, Brian J Haas, Moran Yassour, Joshua Z Levin, Dawn A Thompson, Ido Amit,  
620 Xian Adiconis, et al. 2011. 'Full-Length Transcriptome Assembly from RNA-Seq Data without  
621 a Reference Genome'. *Nature Biotechnology* 29 (7): 644–52.  
622 [http://www.nature.com/nbt/journal/v29/n7/abs/nbt.1883.html#supplementary-](http://www.nature.com/nbt/journal/v29/n7/abs/nbt.1883.html#supplementary-information)  
623 [information](http://www.nature.com/nbt/journal/v29/n7/abs/nbt.1883.html#supplementary-information).
- 624 Guan, Dengfeng, Shane A. McCarthy, Jonathan Wood, Kerstin Howe, Yadong Wang, and Richard  
625 Durbin. 2020. 'Identifying and Removing Haplotypic Duplication in Primary Genome  
626 Assemblies'. *Bioinformatics*, January. <https://doi.org/10.1093/bioinformatics/btaa025>.
- 627 Haas, Brian. 2010. 'TransposonPSI: An Application of PSI-Blast to Mine (Retro-)Transposon ORF  
628 Homologies'. 2010. <http://transposonpsi.sourceforge.net/>.
- 629 Hallab, Asis. 2014. 'Protein Function Prediction Using Phylogenomics, Domain Architecture Analysis,  
630 Data Integration, and Lexical Scoring'. PhD thesis, Bonn.
- 631 Ikegami, F, A Yamamoto, Y H Kuo, and F Lambein. 1999. 'Enzymatic Formation of 2,3-  
632 Diaminopropionic Acid, the Direct Precursor of the Neurotoxin Beta-ODAP, in Lathyrus  
633 Sativus'. *Biological & Pharmaceutical Bulletin* 22 (7): 770–71.
- 634 Ikegami, Fumio, Godelieve Ongena, Ritsuko Sakai, Satoshi Itagaki, Masuko Kobori, Tsutomu Ishikawa,  
635 Yu-Haey Kuo, Fernand Lambein, and Isamu Murakoshi. 1993. 'Biosynthesis of  $\beta$ -(Isoxazolin-  
636 5-on-2-YI)-l-Alanine by Cysteine Synthase in Lathyrus Sativus'. *Phytochemistry* 33 (1): 93–98.
- 637 Jiao, C.-J. J, J.-L. L Jiang, L.-M. M Ke, W. Cheng, F.-M. M Li, Z.-X. X Li, and C.-Y. Y Wang. 2011. 'Factors  
638 Affecting Beta-ODAP Content in Lathyrus Sativus and Their Possible Physiological  
639 Mechanisms'. *Food and Chemical Toxicology* 49 (3): 543–49.  
640 <https://doi.org/10.1016/j.fct.2010.04.050>.
- 641 Jones, Philip, David Binns, Hsin-Yu Chang, Matthew Fraser, Weizhong Li, Craig McAnulla, Hamish  
642 McWilliam, et al. 2014. 'InterProScan 5: Genome-Scale Protein Function Classification'.  
643 *Bioinformatics* 30 (9): 1236–40. <https://doi.org/10.1093/bioinformatics/btu031>.
- 644 Kim, Daehwan, Ben Langmead, and Steven L Salzberg. 2015. 'HISAT: A Fast Spliced Aligner with Low  
645 Memory Requirements'. *Nature Methods* 12 (4): 357–60.  
646 <https://doi.org/10.1038/nmeth.3317>.
- 647 Kislev, Mordechai E. 1989. 'Origins of the Cultivation of Lathyrus Sativus and L. Cicera (Fabaceae)'.  
648 *Economic Botany* 43 (2): 262–70.
- 649 Kolmogorov, Mikhail, Jeffrey Yuan, Yu Lin, and Pavel A. Pevzner. 2019. 'Assembly of Long, Error-  
650 Prone Reads Using Repeat Graphs'. *Nature Biotechnology* 37 (5): 540–46.  
651 <https://doi.org/10.1038/s41587-019-0072-8>.
- 652 Kopylova, Evguenia, Laurent Noé, and Hélène Touzet. 2012. 'SortMeRNA: Fast and Accurate Filtering  
653 of Ribosomal RNAs in Metatranscriptomic Data'. *Bioinformatics* 28 (24): 3211–17.  
654 <https://doi.org/10.1093/bioinformatics/bts611>.
- 655 Kreplak, Jonathan, Mohammed-Amin Madoui, Petr Cápál, Petr Novák, Karine Labadie, Grégoire  
656 Aubert, Philipp E. Bayer, et al. 2019. 'A Reference Genome for Pea Provides Insight into  
657 Legume Genome Evolution'. *Nature Genetics* 51 (9): 1411–22.  
658 <https://doi.org/10.1038/s41588-019-0480-1>.
- 659 Kumar, Shiv, G. Bejiga, S. Ahmed, H. Nakkoul, and A. Sarker. 2011. 'Genetic Improvement of Grass  
660 Pea for Low Neurotoxin ( $\beta$ -ODAP) Content'. *Food and Chemical Toxicology* 49 (3): 589–600.  
661 <https://doi.org/10.1016/J.FCT.2010.06.051>.
- 662 Kuo, Yu-Haey, and Fernand Lambein. 1991. 'Biosynthesis of the Neurotoxin  $\beta$ -N-Oxalyl- $\alpha$ ,  $\beta$ -  
663 Diaminopropionic Acid in Callus Tissue of Lathyrus Sativus'. *Phytochemistry* 30 (10): 3241–  
664 44.
- 665 Kusama-Eguchi, Kuniko, Takaaki Miyano, Makoto Yamamoto, Atsuhiko Suda, Yoshihisa Ito, Kumiko  
666 Ishige, Mayuko Ishii, et al. 2014. 'New Insights into the Mechanism of Neurolathyrism: L- $\beta$ -

- 667 ODAP Triggers [Ca<sup>2+</sup>]accumulation and Cell Death in Primary Motor Neurons through  
668 Transient Receptor Potential Channels and Metabotropic Glutamate Receptors'. *Food and*  
669 *Chemical Toxicology* 67: 113–22. <https://doi.org/10.1016/j.fct.2014.02.021>.
- 670 Lambein, F, J K Khan, Y H Kuo, C G Campbell, and C J Briggs. 1993. 'Toxins in the Seedlings of Some  
671 Varieties of Grass Pea (*Lathyrus Sativus*)'. *Nat Toxins* 1 (4): 246–49.
- 672 Leitch, I.J., E. Johnston, J. Pellicer, O. Hidalgo, and Bennett M.D. 2019. 'Plant DNA C-Values Database,  
673 Release 7.1'. April 2019. <https://data.kew.org/cvalues/>.
- 674 Letunic, Ivica, and Peer Bork. 2019. 'Interactive Tree Of Life (ITOL) v4: Recent Updates and New  
675 Developments'. *Nucleic Acids Research* 47 (W1): W256–59.  
676 <https://doi.org/10.1093/nar/gkz239>.
- 677 Li, Heng. 2013. 'Aligning Sequence Reads, Clone Sequences and Assembly Contigs with BWA-MEM'.  
678 *ArXiv:1303.3997 [q-Bio]*, May. <http://arxiv.org/abs/1303.3997>.
- 679 ———. 2018. 'Minimap2: Pairwise Alignment for Nucleotide Sequences'. *Bioinformatics* 34 (18):  
680 3094–3100. <https://doi.org/10.1093/bioinformatics/bty191>.
- 681 Love, R. Rebecca, Neil I. Weisenfeld, David B. Jaffe, Nora J. Besansky, and Daniel E. Neafsey. 2016.  
682 'Evaluation of DISCOVAR de Novo Using a Mosquito Sample for Cost-Effective Short-Read  
683 Genome Assembly'. *BMC Genomics* 17 (1): 187. <https://doi.org/10.1186/s12864-016-2531-7>.
- 684 Lowe, T M, and S R Eddy. 1997. 'TRNAscan-SE: A Program for Improved Detection of Transfer RNA  
685 Genes in Genomic Sequence.' *Nucleic Acids Research* 25 (5): 955–64.
- 686 Macas, Jiří, and Pavel Neumann. 2007. 'Ogre Elements — A Distinct Group of Plant Ty3/Gypsy-like  
687 Retrotransposons'. *Gene, ASILOMAR* 2006, 390 (1): 108–16.  
688 <https://doi.org/10.1016/j.gene.2006.08.007>.
- 689 Macas, Jiří, Petr Novák, Jaume Pellicer, Jana Čížková, Andrea Koblížková, Pavel Neumann, Iva Fuková,  
690 Jaroslav Doležel, Laura J. Kelly, and Ilia J. Leitch. 2015. 'In Depth Characterization of  
691 Repetitive DNA in 23 Plant Genomes Reveals Sources of Genome Size Variation in the  
692 Legume Tribe Fabeae'. Edited by Andreas Houben. *PLOS ONE* 10 (11): e0143424.  
693 <https://doi.org/10.1371/journal.pone.0143424>.
- 694 Malathi, K, Padmanab.G, and P S Sarma. 1970. 'Biosynthesis of Beta-N-Oxalyl-L-Alpha,Beta-  
695 Diaminopropionic Acid, *Lathyrus Sativus* Neurotoxin'. *Phytochemistry* 9 (7): 1603–9.  
696 [https://doi.org/Doi.10.1016/S0031-9422\(00\)85283-8](https://doi.org/Doi.10.1016/S0031-9422(00)85283-8).
- 697 Mapleson, Daniel, Gonzalo Garcia Accinelli, George Kettleborough, Jonathan Wright, and Bernardo J.  
698 Clavijo. 2017. 'KAT: A K-Mer Analysis Toolkit to Quality Control NGS Datasets and Genome  
699 Assemblies'. *Bioinformatics* 33 (4): 574–76. <https://doi.org/10.1093/bioinformatics/btw663>.
- 700 Mapleson, Daniel, Luca Venturini, Gemy Kaithakottil, and David Swarbreck. 2018. 'Efficient and  
701 Accurate Detection of Splice Junctions from RNA-Seq with Portcullis'. *GigaScience* 7 (12).  
702 <https://doi.org/10.1093/gigascience/giy131>.
- 703 Nandini, A V, B G Murray, I E W O'Brien, and K R W Hammett. 1997. 'Intra- and Interspecific  
704 Variation in Genome Size in *Lathyrus* (Leguminosae)'. *Botanical Journal of the Linnean*  
705 *Society* 125 (4): 359–66. <https://doi.org/10.1111/j.1095-8339.1997.tb02265.x>.
- 706 Nazar, Igor, Maria Claudia, F Emer, Silvia Vergilio, and Mario Jino. 2020. 'XTool: Uma Ferramenta de  
707 Teste Baseado Em Defeitos Para Esquemas de Dados', April.
- 708 Neumann, Pavel, Andrea Koblížková, Alice Navrátilová, and Jiří Macas. 2006. 'Significant Expansion of  
709 *Vicia Pannonica* Genome Size Mediated by Amplification of a Single Type of Giant  
710 Retroelement'. *Genetics* 173 (2): 1047–56. <https://doi.org/10.1534/genetics.106.056259>.
- 711 Neumann, Pavel, Dana Požárková, and Jiří Macas. 2003. 'Highly Abundant Pea LTR Retrotransposon  
712 Ogre Is Constitutively Transcribed and Partially Spliced'. *Plant Molecular Biology* 53 (3): 399–  
713 410. <https://doi.org/10.1023/B:PLAN.0000006945.77043.ce>.
- 714 Ochatt, S J, C Conreux, and L Jacas. 2013. 'Flow Cytometry Distinction between Species and between  
715 Landraces within *Lathyrus* Species and Assessment of True-to-Typeness of in Vitro  
716 Regenerants'. *Plant Systematics and Evolution* 299 (1): 75–85.  
717 <https://doi.org/10.1007/s00606-012-0704-7>.

- 718 Perteza, Mihaela, Geo M. Perteza, Corina M. Antonescu, Tsung-Cheng Chang, Joshua T. Mendell, and  
719 Steven L. Salzberg. 2015. 'StringTie Enables Improved Reconstruction of a Transcriptome  
720 from RNA-Seq Reads'. *Nature Biotechnology* 33 (3): 290–95.  
721 <https://doi.org/10.1038/nbt.3122>.
- 722 Roberts, Adam, Harold Pimentel, Cole Trapnell, and Lior Pachter. 2011. 'Identification of Novel  
723 Transcripts in Annotated Genomes Using RNA-Seq'. *Bioinformatics* 27 (17): 2325–29.  
724 <https://doi.org/10.1093/bioinformatics/btr355>.
- 725 Ruan, Jue, and Heng Li. 2020. 'Fast and Accurate Long-Read Assembly with Wtdbg2'. *Nature*  
726 *Methods* 17 (2): 155–58. <https://doi.org/10.1038/s41592-019-0669-3>.
- 727 Santha, I M, and S L Mehta. 2001. 'Development of Low ODAP Somaclones of Lathyrus Sativus'.  
728 *Lathyrus Lathyrism Newsletter*, no. 2: 42.
- 729 Sarkar, Abhimanyu, Peter M. F. Emmrich, Ashutosh Sarker, Xuxiao Zong, Cathie Martin, and Trevor L.  
730 Wang. 2019. 'Grass Pea: Remodeling an Ancient Insurance Crop for Climate Resilience'. In  
731 *Genomic Designing of Climate-Smart Pulse Crops*, edited by Chittaranjan Kole, 425–69.  
732 Cham: Springer International Publishing. [https://doi.org/10.1007/978-3-319-96932-9\\_9](https://doi.org/10.1007/978-3-319-96932-9_9).
- 733 Sawant, P V, V S Jayade, and S R Patil. 2011. 'Line × Tester Analysis in Lathyrus'. *Journal of Food*  
734 *Legumes* 24 (1): 41–45.
- 735 Sayers, Eric W., Richa Agarwala, Evan E. Bolton, J. Rodney Brister, Kathi Canese, Karen Clark, Ryan  
736 Connor, et al. 2019. 'Database Resources of the National Center for Biotechnology  
737 Information'. *Nucleic Acids Research* 47 (D1): D23–28. <https://doi.org/10.1093/nar/gky1069>.
- 738 Seoane, Pedro, Noe Fernandez, and Dario Guerrero. 2018. *Full\_length\_next* (version 1.0.1).  
739 [https://rubygems.org/gems/full\\_length\\_next/versions/1.0.1](https://rubygems.org/gems/full_length_next/versions/1.0.1).
- 740 Shafin, Kishwar, Trevor Pesout, Ryan Lorig-Roach, Marina Haukness, Hugh E. Olsen, Colleen  
741 Bosworth, Joel Armstrong, et al. 2019. 'Efficient de Novo Assembly of Eleven Human  
742 Genomes Using PromethION Sequencing and a Novel Nanopore Toolkit'. *BioRxiv*, July,  
743 715722. <https://doi.org/10.1101/715722>.
- 744 Siddique, K H M, C L Hanbury, and A Sarker. 2006. 'Registration of "Ceora" Grass Pea Registration by  
745 CSSA'. *Crop Science* 46 (2): 986. <https://doi.org/10.2135/cropsci2005.0131>.
- 746 Silvestre, Susana, Susana de Sousa Araújo, Maria Carlota Vaz Patto, and Jorge Marques da Silva.  
747 2014. 'Performance Index: An Expedient Tool to Screen for Improved Drought Resistance in  
748 the Lathyrus Genus'. *Journal of Integrative Plant Biology* 56 (7): 610–21.
- 749 Simão, Felipe A, Robert M Waterhouse, Panagiotis Ioannidis, Evgenia V Kriventseva, and Evgeny M  
750 Zdobnov. 2015. 'BUSCO: Assessing Genome Assembly and Annotation Completeness with  
751 Single-Copy Orthologs'. *Bioinformatics* 31 (19): 3210–12.  
752 <https://doi.org/10.1093/bioinformatics/btv351>.
- 753 Simpson, Jared T., Kim Wong, Shaun D. Jackman, Jacqueline E. Schein, Steven J.M. Jones, and Inanç  
754 Birol. 2009. 'ABYSS: A Parallel Assembler for Short Read Sequence Data'. *Genome Research*  
755 19 (6): 1117–23. <https://doi.org/10.1101/gr.089532.108>.
- 756 Slater, Guy St C, and Ewan Birney. 2005. 'Automated Generation of Heuristics for Biological  
757 Sequence Comparison'. *BMC Bioinformatics* 6 (February): 31. <https://doi.org/10.1186/1471-2105-6-31>.
- 759 Smit, A.F.A., R. Hubley, and P. Green. 2013. *RepeatMasker* (version open 4.0.0).
- 760 Stanke, Mario, Oliver Keller, Irfan Gunduz, Alec Hayes, Stephan Waack, and Burkhard Morgenstern.  
761 2006. 'AUGUSTUS: Ab Initio Prediction of Alternative Transcripts'. *Nucleic Acids Research* 34  
762 (Web Server issue): W435–39. <https://doi.org/10.1093/nar/gkl200>.
- 763 Stanke, Mario, Rasmus Steinkamp, Stephan Waack, and Burkhard Morgenstern. 2004. 'AUGUSTUS: A  
764 Web Server for Gene Finding in Eukaryotes'. *Nucleic Acids Research* 32 (suppl 2): W309–12.  
765 <https://doi.org/10.1093/nar/gkh379>.
- 766 Sukumaran, Jeet, and Mark T. Holder. 2010. 'DendroPy: A Python Library for Phylogenetic  
767 Computing'. *Bioinformatics (Oxford, England)* 26 (12): 1569–71.  
768 <https://doi.org/10.1093/bioinformatics/btq228>.



- 769 The UniProt Consortium. 2019. 'UniProt: A Worldwide Hub of Protein Knowledge'. *Nucleic Acids*  
770 *Research* 47 (D1): D506–15. <https://doi.org/10.1093/nar/gky1049>.
- 771 Trapnell, Cole, Brian A. Williams, Geo Pertea, Ali Mortazavi, Gordon Kwan, Marijke J. van Baren,  
772 Steven L. Salzberg, Barbara J. Wold, and Lior Pachter. 2010. 'Transcript Assembly and  
773 Quantification by RNA-Seq Reveals Unannotated Transcripts and Isoform Switching during  
774 Cell Differentiation'. *Nature Biotechnology* 28 (5): 511–15.  
775 <https://doi.org/10.1038/nbt.1621>.
- 776 Tsegaye, D, W Tadesse, and M Bayable. 2005. 'Performance of Grass Pea (*Lathyrus Sativus* L.)  
777 Somaclones at Adet, Northwest Ethiopia'. *Lathyrus Lathyrism Newsletter* 4: 5–6.
- 778 Vaz Patto, Maria Carlota, M Fernández-Aparicio, A Moral, and D Rubiales. 2006. 'Characterization of  
779 Resistance to Powdery Mildew (*Erysiphe Pisi*) in a Germplasm Collection of *Lathyrus Sativus*'.  
780 *Plant Breeding* 125 (3): 308–10.
- 781 Venturini, Luca, Shabhonam Caim, Gemy George Kaithakottil, Daniel Lee Mapleson, and David  
782 Swarbreck. 2018. 'Leveraging Multiple Transcriptome Assembly Methods for Improved Gene  
783 Structure Annotation'. *GigaScience* 7 (8). <https://doi.org/10.1093/gigascience/giy093>.
- 784 Venturini, Luca, Gemy Kaithakottil, and David Swarbreck. 2016. 'Extended Methods for the  
785 Annotation of *Triticum Aestivum* CS42.' Earlham Insitute, Norwich, UK.
- 786 Vondrak, Tihana, Laura Ávila Robledillo, Petr Novák, Andrea Koblížková, Pavel Neumann, and Jiří  
787 Macas. 2020. 'Characterization of Repeat Arrays in Ultra-Long Nanopore Reads Reveals  
788 Frequent Origin of Satellite DNA from Retrotransposon-Derived Tandem Repeats'. *The Plant*  
789 *Journal* 101 (2): 484–500. <https://doi.org/10.1111/tpj.14546>.
- 790 Wysokar, A., K. Tibbetts, M. McCown, N. Homer, and T. Fennell. (2016) 2016. *Picard: A Set of Java*  
791 *Command Line Tools for Manipulating High-Throughput Sequencing Data (HTS) Data and*  
792 *Formats*. Java. <https://github.com/kcibul/picard>.
- 793 Yadav, S S, G Bejiga, M Brink, and G Belay. 2006. 'Lathyrus Sativus L.' PROTA4U. Wageningen,  
794 Netherlands: PROTA (Plant Resources of Tropical Africa / Ressources végétales de l'Afrique  
795 tropicale). 2006. <http://www.prota4u.org/search.asp>.
- 796 Yang, Hui-Min, and Xiao-Yan Zhang. 2005. 'Considerations on the Reintroduction of Grass Pea in  
797 China'. *Lathyrus Lathyrism Newsletter* 4: 22–26.
- 798 Zhelyazkova, Tsenka, Dimitar Pavlov, Grosi Delchev, and Antoniya Stoyanova. 2016. 'Productivity and  
799 Yield Stability of Six Grain Legumes in the Moderateclimatic Conditions in Bulgaria'. *Scientific*  
800 *Papers. Series A. Agronomy* LIX: 478–87.
- 801 Zimin, Aleksey V., Guillaume Marçais, Daniela Puiu, Michael Roberts, Steven L. Salzberg, and James  
802 A. Yorke. 2013. 'The MaSuRCA Genome Assembler'. *Bioinformatics* 29 (21): 2669–77.  
803 <https://doi.org/10.1093/bioinformatics/btt476>.
- 804

805

806