

1 **Transcriptome sequencing and marker development in winged bean (*Psophocarpus***
2 ***tetragonolobus*; Leguminosae)**

3
4 Mohammad Vatanparast¹, Prateek Shetty², Ratan Chopra³, Jeff J. Doyle⁴, N.
5 Sathyanarayana^{5*} and Ashley N. Egan^{1*}

6
7 ¹US National Herbarium (US), Department of Botany, Smithsonian Institution-NMNH, 10th
8 and Constitution Ave, Washington DC, 20013, USA.

9 ²Department of Plant Biology, Michigan State University, 612 Wilson Road, Room 166, East
10 Lansing, MI, 48824, USA.

11 ³United States Department of Agriculture, Agriculture Research Service, 3810 4th St.,
12 Lubbock, TX, 79415, USA.

13 ⁴Section of Plant Breeding & Genetics, School of Integrative Plant Science, Cornell
14 University, 412 Mann Library, Ithaca, NY, 14853, USA.

15 ⁵Department of Botany, Sikkim University, 5th Mile, Tadong, Gangtok, Sikkim, 737102,
16 India.

17
18 *Email Addresses:*

19 Mohammad Vatanparast: Vatanparastm@si.edu

20 Prateek Shetty: prateekshettys@gmail.com

21 Ratan Chopra: Ratan.Chopra@ARS.USDA.GOV

22 Jeff J. Doyle: jjd5@cornell.edu

23 N. Sathyanarayana: nsathyanarayana@cus.ac.in

24 Ashley N. Egan: egana@si.edu

25

26 **Corresponding authors:*

27 Ashley N. Egan (primary)

28 N. Sathyanarayana (secondary)

29 Winged bean, *Psophocarpus tetragonolobus* (L.) DC, is similar to soybean in yield and
30 nutritional value but more viable in tropical conditions. Here, we strengthen genetic resources
31 for this orphan crop by producing a *de novo* transcriptome assembly and annotation of two
32 Sri Lankan accessions (denoted herein as CPP34 [PI 491423] and CPP37 [PI 639033]),
33 developing simple sequence repeat (SSR) markers, and identifying single nucleotide
34 polymorphisms (SNPs) between geographically separated genotypes. A combined assembly
35 based on 804,757 reads from two accessions produced 16,115 contigs with an N50 of 889 bp,
36 over 90% of which has significant sequence similarity to other legumes. Combining contigs
37 with singletons produced 97,241 transcripts. We identified 12,956 SSRs, including 2,594
38 repeats for which primers were designed and 5,190 high-confidence SNPs between Sri
39 Lankan and Nigerian genotypes. The transcriptomic data sets generated here provide new
40 resources for gene discovery and marker development in this orphan crop, and will be vital
41 for future plant breeding efforts. We also analyzed the soybean trypsin inhibitor (STI) gene
42 family, important plant defense genes, in the context of related legumes and found evidence
43 for radiation of the Kunitz trypsin inhibitor (KTI) gene family within winged bean.

44

45 Winged bean (*Psophocarpus tetragonolobus* (L.) DC) is a promising legume crop of the
46 world's tropical regions. It is predominantly self-pollinated and possesses a twining habit,
47 tuberous roots, longitudinally winged pods, and both annual and perennial growth forms ¹.
48 The genus *Psophocarpus* Neck. ex DC comprises 10 species. Excluding cultivated winged
49 bean, all other species are wild and native to Africa, Madagascar and the Mascarene Islands
50 in the Indian Ocean ². Winged bean is speculated to have originated from the progenitor
51 species *P. grandiflorus* R. Wilczek and is now cultivated extensively in Papua New Guinea
52 and Southeast Asia, and to a lesser extent in Africa ^{1,2}. Winged bean has a diploid genome
53 ($2n=2x=18$) ³ and an estimated genome size of 1.22 Gbp/C (A.N. Egan, unpublished data).

54

55 Every part of the winged bean is edible, earning it the distinction of “*Supermarket on a stalk*”
56 ⁴. The exceptional nutritional quality of this plant, and the fact that it provides suitable human
57 food sources at all stages of its life cycle, makes it a promising candidate for increased,
58 widespread use in protein deficient tropical areas of the world. The young pods contain 2.4
59 grams (g) protein per 100g of edible portion; the dried tubers and seeds contain 8-20% and
60 34% protein, respectively, as well as a high oil contents (18%) - traits which have earned it
61 the name “*soybean of the tropics*” ⁵. If both seed and tuber yields are combined, winged bean
62 can outperform many other legume crops that are conventionally grown in the tropics and
63 thus offers a cheap nutritional food source. Consequently, it is projected as a promising
64 alternative to soybean in areas where soybean cultivation success is marginal.

65

66 Since the 1975 publication by the US National Academy of Sciences of *The Winged Bean: A*
67 *High Protein Crop for the Tropics* ⁶, considerable effort has been focused on studying the
68 nutritional quality and climatic and ecological tolerances of the plant ^{7,8}. Winged bean
69 reportedly possesses anti-nutritional factors such as phytolectinins, cyanogenic glycosides,

70 tannins, lectins, flatulence factors, and saponins ⁹. However, processing using moist heat or
71 soaking has been shown to safely eliminate these substances . Research efforts concerning
72 such anti-nutritional components have yielded significant knowledge concerning trypsin, a
73 serine protease that acts to hydrolyze proteins as part of vertebrate digestion, and trypsin
74 inhibitors, proteins that stop the action of trypsin, thereby interfering with digestion. It has
75 been suggested that trypsin inhibitors play a role in protecting plant tissues against the action
76 of bacterial proteases at the colonization site of pathogenic bacteria ¹⁰. In addition, studies
77 show involvement in defense against insects that suck the phloem sap and against bacteria
78 that invade upon wounding ¹¹. In biomedical research, these modes of action have made
79 trypsin and trypsin inhibitors vital components of molecular cell research where they are
80 widely used in cell culture to detach cells from tissue culture plates. Since their first
81 discovery in soybeans in 1945 ¹², other Kunitz-type trypsin inhibitors have been discovered
82 and characterized from winged bean ^{13,14}, predominantly from seeds.

83

84 It is hard to find another high rainfall-adapted tropical legume with as many desirable
85 characteristics as winged bean¹. However, much needs to be done in terms of breeding
86 efforts, especially to develop self-supporting, determinate cultivars bearing large numbers of
87 relatively small pods having nutritious seeds and tubers, and cultivars resistant to biotic and
88 abiotic stresses. Considerable variability for growth vigor and quantitative characters such as
89 protein and oil content as well as photoperiodic responses has been recorded ¹⁵. Several
90 beneficial mutants were recovered during the 1970s and '80s through mutation breeding
91 experiments ¹⁶. However, a recent study using inter-simple sequence repeat (ISSR) markers
92 reported low genetic diversity among the winged bean germplasm collected from different
93 parts of the world ¹⁷. With the advent of genomic tools such as molecular markers, genetic
94 maps etc., the genetic improvement of underutilized crops has been greatly facilitated,

95 enabling the development of improved genotypes or varieties with enhanced trait values¹⁸. In
96 the case of winged bean, studies on genomic resource development for enabling basic and
97 applied research on genetics, evolution, ecology and molecular breeding programs are
98 lacking, yet the advent of genomic technologies provides significant prospects for
99 improvement¹⁹. Transcriptome sequencing is cost-effective and a valuable method for
100 efficient and rapid identification of molecular markers in resource poor plant species²⁰.

101

102 The present study was undertaken with the following objectives: (a) to generate a set of
103 expressed sequence tag (EST) resources through whole transcriptome analysis based on
104 Roche 454-based transcriptomes for two winged bean accessions from Sri Lanka; (b) to
105 develop a *de novo* assembly for these transcriptomes; (c) to annotate the transcriptome
106 information; and (d) to discover microsatellite markers for future genetic studies. We also
107 compared Sri Lankan accessions to a Nigerian winged bean transcriptome previously
108 sequenced on the Illumina platform (e) to identify Single Nucleotide Polymorphisms (SNPs)
109 evident between the geographically separated genotypes and (f) to present an analysis of the
110 Kunitz trypsin inhibitor gene family in the context of related legumes.

111

112 **Results**

113

114 *Sequencing and De novo assembly of winged bean transcriptomes*

115

116 Pyrosequencing of two Sri Lankan accessions produced comparable sequence output, where
117 genotype CPP34 produced a total of 369,820 single-end reads comprising 136,943,216 bp
118 with an average read length of 574 bp and genotype CPP37 produced a total of 334,639
119 single-end reads comprising 92,126,948 bp with an average read length of 565 bp (Table 1).

120 Using read count as a proxy, the depth of sequencing across our contigs was similar for the
121 independent *de novo* assemblies, ranging from one to 4,953 reads, with an average read depth
122 of 25 reads per contig for CPP34 and ranging from one to 3,972 reads with an average read
123 depth of 30 reads per contig for CPP37. Comparison of transcripts from the CPP34 and
124 CPP37 independent assemblies (Supplementary file 1, inclusive of Tables S1-S3 and Figs.
125 S1, online) found fewer than 200 high-confidence SNPs between them (data not shown),
126 equating to approximately one SNP every 150,000 bp. Therefore, reads from the
127 independently sequenced accessions were combined and co-assembled. For the combined
128 assembly (CPP34-7), this translated to 704,459 reads comprising 229,070,164 bases from
129 both accessions (Table 1). Because 454 pyrosequencing produces comparatively long reads
130 (300-800 bp long), unassembled reads, here notated as singletons post-assembly, may
131 potentially represent full-length mRNA transcripts. In order to not lose potential information,
132 singletons of the CPP34-7 were extracted and appended to the final assembly of CPP34-7 and
133 used in the GO and SNP analyses.

134

135 *Functional annotation & legume sequence similarity*

136

137 For the GO analysis, the combined assembly of CPP34-7 was used with inclusion of
138 singletons (16,115 contigs plus 81,126 singletons, table 1). Using a total of 97,241
139 transcripts, TransDecoder could track 33,038 transcripts against BLAST and Pfam databases.
140 Of these 33,038 transcripts, BLAST searches against NCBI's nr database retrieved 32,993
141 transcripts with hits (see Supplementary file 2 online), discarding 45 transcripts that had zero
142 hits in NCBI. Therefore, 64,248 (66%) of our original 97,241 transcripts did not hit any
143 known gene or DNA region in NCBI and Pfam databases, of which 62,783 were singletons.
144 Thus, 79% of singletons were discarded in the BLAST searching steps due to a lack of

145 annotation. Of the 32,993 transcripts with BLAST hits, the GO analysis determined GO ID
146 and enzyme code (EC) assignments for 16,561 (50.1%) with full or partial annotations (Fig. 1
147 in text, and see Supplementary file 2 online). Of the 16,561 annotated transcripts, 5,053 have
148 predicted functions (EC codes). Overall, 2,829 transcripts were not functionally annotated by
149 Blast2GO (zero hits) of which 1,932 (68%) corresponded to singletons. Participation of genes
150 in a particular biological process and molecular function are shown in figure 2. Several
151 transcripts were assigned to more than one GO term; therefore, the total number of GO terms
152 obtained for our dataset was higher than the total number of transcripts. In total, 47,178 GO
153 terms were retrieved, with 46.2%, 37% and 16.8%, corresponding to the MF, BP, and CC
154 categories, respectively. In the MF category, nucleotide binding (number of sequences =
155 3,413), kinase activity (1,474) and DNA binding (1,200) had the highest number of assigned
156 sequences. In the BP category, cellular protein modification (1,953), carbohydrate metabolic
157 processes (1,080) and transport (908) were the majority and in the case of CC, genes involved
158 in the plastid (319), cytoskeleton (288) and ribosome (281) activities were highly represented
159 (Fig. 2).

160

161 A comparison of our assembled contigs against other legume NCBI protein sequence
162 databases from chickpea, pigeon pea, soybean, common bean, *Medicago truncatula*, and
163 *Lotus japonicus* using the BLASTX program from NCBI showed that 15,558 of 16,115
164 (96.5%) contigs from the CPP34-7 assembly had significant sequence similarity to sequences
165 in one or more legume protein databases. About 90.5% of the 16,115 contigs had $\geq 80\%$
166 sequence identity (Fig. 3). The majority of the contigs (57.3%) were most similar to *Glycine*
167 *max* (Fig. 4), a finding that, at first glance, seems to contradict that expected based on
168 evolutionary relationships of legume lineages, but is likely due to the relative over-
169 representation of genes within the soybean genome due to i) recent whole genome

170 duplication and ii) a much higher level and standard of annotation and gene discovery
171 relative to other legume genomes. Differences in evolutionary rate across lineages may also
172 impact this outcome. In relation to *Phaseolus vulgaris*, it is known that *Phaseolus* has a
173 higher mutation rate than *Glycine* and related lineages,^{21,22} which could increase the
174 divergence, and thereby decrease the best-BLAST hits, of *Psophocarpus* against *Phaseolus*
175 relative to *Glycine*. However, this explanation is invoked with caution given that it assumes
176 similar relative rates between *Glycine* and *Psophocarpus*, information that is beyond the
177 scope of this project.

178

179 *Identification of transcription factors*

180

181 In the overall GO analysis, 274 transcripts were annotated as transcription factors (Fig. 2). Of
182 the 16,115 contigs, 176 putative winged bean transcription factor genes, distributed in at least
183 ten families, were identified representing 1.1% of winged bean transcripts, which were
184 assigned to different categories. Among these, basic leucine zipper (bZIP; 32), Teosinte-
185 Branched1/Cycloidea/PCF (TCP; 19), MADS (17), MYB (11) and WRKY (9) were among
186 the top five categories (Fig. 5.). The overall distribution of transcription factor encoding
187 transcripts among the various known protein families is very similar to that of soybean and
188 other legumes. However, almost all families showed minor species specific differences (for
189 example, bZIP, MYB, WRKY etc.) with regard to TF gene families reported for *Lotus*,
190 *Medicago* and *Glycine max*.

191

192 *Identification of simple sequence repeats*

193

194 The SSR analysis detected 10,984 perfect SSRs, 13 imperfect SSRs, and 1,959 compound

195 SSRs, for a total of 12,956 SSRs (see Supplementary file 3 online). Of the 10,984 primary
196 SSRs, 57 were adenine (A: 30) or thymine (T: 27) monomers with at least 13 repeats. These
197 were assumed to represent remnants of mRNA poly-A tails and were thus removed prior to
198 primer prediction. No runs of 12 or more cytosine or guanine monomeric repeats were found.
199 Nearly three-quarters of the remaining 10,927 perfect SSRs (7,933) were hexamers with only
200 two repeats. Although these 12-mers may be useful as linkage markers, the low number of
201 repeat units would likely take these out of the microsatellite category. The remaining 2,994
202 perfect SSRs were distributed across di-, tri-, tetra-, penta-, and hexamer SSRs (Fig. 6) and
203 were used for primer creation. The majority (63%) of SSRs were detected in the tri- and
204 hexamer categories (Fig 6a). In general, the number of SSRs detected in each size category
205 decreased with increasing repeat number (Figure 6b-f). Primers were successfully created for
206 2,594 SSRs with product sizes ranging from 100 to 280 bp (see Supplementary file 4 online).
207 Analysis of the primed SSRs showed bias towards certain di- and tri- repeat type motifs
208 (Table 2).

209

210 *Single nucleotide polymorphism discovery*

211

212 GS Reference Mapper mapped 87.7% of reads from Chapman²³ onto the CPP34-7 reference
213 ‘genome’ which consisted of the 97,241 transcripts. Of the 14,571,393 bp of mapped reads,
214 we identified 113,757 SNPs with >95% confidence from the 454HCDiffs file (available upon
215 request), suggesting a SNP frequency of one in 128 bp of coding regions. Interestingly, the
216 majority of high-confidence SNPs were found within singletons (91,686; 80.6%) vs. contigs
217 (22,071; 19.4%), a higher percentage than expected given that singletons make up 67.9% of
218 total transcript length. As a conservative measure, we filtered SNPs based on allele frequency
219 from >95% - 100% confidence levels and those having >20x coverage (Table 3), producing a

220 total of 13,091 SNPs distributed across 10,176 transcripts of which 5,196 (39.7%) were from
221 contigs, representing 1 SNP every 1,113 bp. The subsequent increase in the proportion of
222 SNPs within contigs is expected in this case given that more highly expressed genes will be
223 more likely to be represented by >20x coverage and are most likely to assemble into contigs.
224 Lastly, we removed all single nucleotide indels (7,665 of the 13,091) and those length
225 variants that involved insertions or deletions of one or more nucleotides alone (i.e. those
226 without point mutations involved in the length variants), resulting in a high-confidence set of
227 5,190 SNPs, 96% of which are one-to-one point mutations (see Supplemental file 5 online).
228 Within the 5,190 SNPs, 151 unique SNP patterns were found and 211 (4%) SNPs were length
229 variants involving one or more point mutation within the length variant. Of the 4,979 one-to-
230 one polymorphisms, 3,433 (68.9%) were transitions and 1,546 (31.1%) were transversions,
231 producing a transition:transversion ratio of 2.22.

232

233 *Kunitz-type trypsin inhibitor gene family analysis*

234

235 We identified 28 contigs from CPP34-7 and 20 contigs from the Chapman²³ transcriptome
236 assembly corresponding to the KTI gene family within the *Psophocarpus* transcriptome (see
237 Supplementary file 6 online). Due to the large number of paralogues in each species, there is
238 no obvious criterion available for rooting this tree, so it was rooted with the largest clade of
239 non-legume sequences, a clade of *Arabidopsis* sequences. Given this rooting, the Bayesian
240 STI gene tree has a polytomous backbone and suggests six distinct subclades based on
241 relatively high posterior probability and bootstrap support, here labeled as A-F (Fig. 7 in text;
242 and see Supplementary file 7 online). The dominant feature of the tree (regardless of rooting)
243 is a lack of clear orthologous relationships across taxa, with evidence of lineage-specific
244 amplification of STI and KTI genes in each species. For example, subclade A comprises two

245 clades made up of only *Populus* sequences and an *Arabidopsis* STI member as well as a
246 number of clades containing *Glycine*, *Phaseolus*, or *Medicago* gene family members, but
247 with no *Psophocarpus* sequences included, whereas subclade C comprises *Populus* sequences
248 only, illustrating a major intra-specific STI radiation (Fig. 7). The vast majority of
249 *Psophocarpus* sequences cluster in clade F, along with many *Glycine* and a single *Phaseolus*
250 sequence. Of the *Psophocarpus* sequences in subclade F, 15 contigs are paired between
251 CPP34-7 and Chapman, forming sister groups that likely represent the same gene in each
252 transcriptome, whereas 13 are unique (Fig. 7). Subclade F illustrates lineage-specific KTI
253 expansion in both *Psophocarpus* and *Glycine*. All *Psophocarpus* sequences obtained from the
254 Pfam or NCBI databases were nested within subclade F, where the majority of the Pfam
255 sequences appeared as monophyletic clades with a contig each from CPP34-7 and Chapman
256 nested therein (Fig. 7).

257

258 **Discussion**

259

260 The legumes represent the third largest family of the flowering plants, many of which are
261 important sources of food, fodder, oil, fiber and medicines. However, with the exception of
262 common pulse crops such as soybean, common bean, etc., a large number of legumes have
263 remained underutilized due to poorly developed infrastructure, especially for genetic and
264 genomic resources²⁴. The advent of genomic technologies has brightened the prospects for
265 such orphan crops^{20,25,26}, with recent research focusing on lentil (*Lens culinaris* Medik.,²⁷),
266 chickpea²⁸, grass pea (*Lathyrus sativus* L.,²⁹), and a number of *Vigna* species^(30,31), among
267 others. Winged bean represents a promising alternative to protein-rich soybean for tropical
268 regions of the world that house nearly 40% of the world population, of which nearly one third

269 is protein deficient, and many of whom are women and children ³². Genomics assisted
270 breeding and enabling biotechnologies that stem from it offer significant promise for targeted
271 genetic improvement of nutritional and other quality traits in winged bean, thus aiding in the
272 development of a low input, high quality legume-based protein diet for these parts of the
273 world. Our combined assembly presents a genetic resource that can be mined for future
274 genetic improvement and plant breeding initiatives. This paper reports development of
275 genetic resources, including molecular markers, in winged bean, in addition to insights into
276 the divergence of the Kunitz type trypsin inhibitors, which are important anti-nutritive agents
277 in winged bean and other legumes.

278

279 In this study, we were able to annotate 32,993 (34%) transcripts from the winged bean
280 combined assembly (CPP34-7; Fig. 2). Schmutz et al ²¹ annotated 27,197 protein-coding
281 genes and 31,638 protein-coding transcripts from the *Phaseolus vulgaris* genome, suggesting
282 that our annotated gene complement is reasonable, although it is likely that our
283 transcriptomes do not provide a full gene complement due to low sequencing depth. Our level
284 of unannotated transcripts is similar to results reported from other non-model organisms,
285 including chickpea ³³ and field pea ³⁴. These unidentified transcripts are likely due to: 1)
286 correspondence to non-coding regions or pseudogenes, 2) short length of transcripts, or 3)
287 novel coding genes that have yet to be described. Cellular, metabolic and transport processes
288 were among the most highly represented groups in terms of GO analysis, as expected given
289 that flower buds, young leaves and shoots are undergoing rapid growth and extensive
290 metabolic activities.

291

292 Singletons (unassembled reads) in *de novo* transcriptome assemblies stem from such

293 phenomena as differences in assembly algorithms, sequencing errors, artifacts in cDNA
294 library construction, gene expression at low levels, or contamination from other organisms
295 such as bacteria or fungi ³⁵. Assessing the GC content of a transcriptome assembly can aid in
296 checking for possible contamination as different organisms have different genomic GC
297 content. We compared the GC content in our data against related legumes to check for
298 contamination and found no evidence of it (see Supplementary file 1 online). In the GO
299 analysis, ~80% of the singletons were discarded in the BLAST step, while the remaining
300 20% persisted, but only 10% proceeded through the GO annotation. Others have found
301 similar low levels of singleton annotation ³⁶, yet, this low level of singleton annotation has
302 lead many to throw out unassembled reads. However, given the comparative length of 454
303 reads, these could easily represent full-length transcripts. Thus, we included singletons in the
304 GO and SNP analyses to evaluate their potential.

305

306 Because transcription factors play important roles in regulating plant functions, we paid
307 particular attention to their number and distribution within winged bean and in relation to
308 other legumes. Several TF gene families are preserved across different plant genera,
309 indicating conserved gene regulatory machinery in plants, as has been shown in legumes
310 previously ³⁷. In this study, we found that 2.4% of the total transcripts are putative
311 transcription factors according to GO analyses, a percentage much lower than the estimated
312 12% found in soybean (based on ~46,430 protein-coding genes) ³⁸. However, Libault et al ³⁷
313 estimated the number of TF-encoding genes across a number of species and found soybean to
314 have 3-4× the number of TFs relative to *Medicago truncatula* or *Lotus japonicus*, likely due
315 to recent whole genome duplication in soybean. If we compare the estimated number of TFs
316 in *M. truncatula* (1473) ³⁷ against the number of putative protein-coding genes in *M.*
317 *truncatula* (~66,000; phytozome v11; <https://phytozome.jgi.doe.gov/pz/portal.html>), we

318 come out with a much more similar estimate (2.2%). However, this is likely an underestimate
319 and a shifting target as the annotation of *M. truncatula* is ongoing.

320

321 Our overall distribution of transcription factors in winged bean within the known TF families
322 is similar to that in soybean and other legume species, with bZIP, MYB, TCP, and WRKY
323 highly represented (^{28,39}). The TF family most highly-represented in our data was the bZIP
324 family, which includes regulators of many central developmental and physiological processes
325 and abiotic and biotic stress responses ^{40,41}. In addition, elevated levels of expression were
326 also found for TCPs and MYB: TCPs have been characterized in other plant species for their
327 role in growth, development, and sex determination ^{42,43}, whereas the MYB family has been
328 implicated in regulation of disease resistance and water loss regulation via stomatal
329 movement ⁴⁴. However, a significant portion of our transcripts comprised several smaller TF
330 families, here classified under the miscellaneous category for want of detailed
331 characterization. Also, we observed minor species-specific differences in the numbers and
332 proportion of our TFs relative to predicted TFs in *Lotus*, *Medicago* and *Glycine max* ³⁷.
333 Further investigation is thus needed to elucidate the evolutionary and functional significance
334 of these events in winged bean.

335

336 Simple sequence repeat, or microsatellite, markers have long been used for genetic diversity
337 analyses and plant breeding efforts, largely due to their highly polymorphic, co-dominant
338 nature, prevalence throughout the genome, ease of use, and cost-effectiveness ⁴⁵. Because
339 they originate in coding regions, SSRs derived from genes have increased amplification
340 success in related species, are useful for assessing functional diversity and for marker-
341 assisted selection, and can act as anchor markers for evolutionary and comparative mapping
342 studies ^{46,47}. While some research has suggested that SSRs derived from coding regions are

343 less polymorphic than their anonymous counterparts ⁴⁷, numerous population genetic,
344 evolutionary, and plant breeding studies have found them to have adequate, if not higher,
345 levels of polymorphism within legumes ^{48,49}.

346

347 Within our data, we discovered nearly 5,000 perfect or compound genic-SSRs with three or
348 more repeats. After filtering for perfect, simple SSRs, we discovered an unequal distribution
349 across size-classes, with trinucleotide repeats making up the bulk (43%) of filtered SSRs.
350 Given the coding nature of the transcriptome, this finding makes sense as proliferation of tri-
351 nucleotide, in-frame repeats would be more tolerated ⁴⁶. The same trend has been noted in
352 other plants, including for legumes *Medicago truncatula* ⁵⁰ and peanut (*Arachis hypogaea* L.;
353 ⁵¹). Of the 2,594 SSRs for which primers were created, 1,928 (74.3%) were annotated in our
354 Blast2GO analyses, 871 (45.2%) of which are putatively homologous to known proteins
355 while 1,057 (54.8%) were similar to hypothetical, uncharacterized, unknown, or predicted
356 proteins (mostly from *Phaseolus vulgaris* and *Glycine max* genome annotations).

357

358 Certain repeat motif types were more prevalent than others in our set of primed SSRs (Table
359 2), a not-uncommon finding that has been documented in legumes previously ⁵². Zhang et al
360 ⁵³ first documented the bias of microsatellites to AG and AAG motifs in *Arabidopsis*, also
361 noting differences of SSR distributions between 5' and 3' untranslated and coding regions,
362 and correlation between trinucleotide repeat motifs and codon usage. In winged bean, the
363 SSR repeat motif type (AG/GA/TC/CT)_n represented the majority (77.7%) of all dinucleotide
364 repeats, while motif types (AT/TA)_n and (AC/CA/GT/TG)_n comprised 13.7% and 8.6%,
365 respectively. Our distribution and ranking of dinucleotide repeat motifs mirrors that in
366 *Arabidopsis* (Table 2). The bias towards the repeat type containing AG and against that of
367 GC has also been found in other plants, including *Phaseolus* ⁵⁴, *Myrciaria dubia* (Kunth)

368 McVaugh⁵⁵ and across eukaryotes⁵⁶. Past research has suggested that AG motifs are most
369 prevalent in 5' untranslated regions^{52,57} and possibly are involved in transcription and
370 regulation⁵³. As mentioned earlier, 25.7% of primed genic-SSR transcripts are unannotated,
371 some of which may correspond to 5' untranslated regions where AG motif types are more
372 prevalent. The frequency of trinucleotide repeat motif types was biased towards AAG in our
373 set of primed SSRs, with this motif type comprising 29.1% of the 10 trinucleotide types,
374 followed by the ATC motif type, comprising 13.9%. The ranking of these and other motifs
375 closely resembles that of *Arabidopsis* (Table 2), with the first two most prevalent motifs the
376 same⁵³.

377

378 SNPs provide another means of assessing genetic variation and, although less polymorphic
379 than SSRs, are abundant and easily obtained via high-throughput sequencing. For example,
380 Rajesh and Muehlbauer⁵⁸ estimated SNP frequency to be one in 66 bp in coding regions and
381 one in 71 bp in genomic regions of chickpea⁵⁸. In another study, Hyten et al⁵⁹ reported
382 7,000-25,000 predicted SNPs through deep resequencing of soybean by a whole genome
383 sequence approach. In this study, we discovered more than 5,190 high-confidence SNPs
384 between our Sri Lankan samples and the geographically separated Nigerian genotype²³. SNP
385 markers identified in this study can be used in quantitative trait loci (QTL) mapping,
386 generating linkage maps, genotyping and breeding studies. Validation of SNPs determined
387 herein is beyond the scope of this paper, nevertheless, this list presents a significant resource
388 for future work in plant breeding and genetic diversity assessment⁶⁰ and marks the first SNP
389 markers discovered to date in *Psophocarpus*.

390

391 Our high-confidence SNP set included 4,979 one-to-one SNPs (those without length variants
392 and involving changes between a single nucleotide position), equating to a

393 transition:transversion (ts:tv) ratio of 2.22. This bias is commonly observed across SNPs
394 throughout a genome, resulting from rampant methyl cytosine to uracil mutations ⁶¹. Similar
395 ratios were found across SNPs in other legumes ^{62,63}. In total, 44% of SNPs identified were
396 found in singletons, a proportion not unexpected given that 68% of transcript read length is in
397 singletons. But the very fact that alignable and putatively homologous singletons were found
398 across the geographically separated and independently sequenced genotypes provides
399 vindication for their inclusion in transcriptome characterizations, at least for 454 data.
400 However, caution is warranted due to the ‘singleness’ of the unassembled read acting as a
401 reference sequence.

402

403 Trypsin inhibitors play important roles in plant development and defense systems and have
404 been studied from various aspects like biotic stress and wounding. These compounds inhibit
405 activity of proteases and are induced by mechanical wounding in leaves, suggesting a strong
406 role as anti-herbivory agents ⁶⁴. Trypsin inhibitors present in legumes include KTIs, the
407 Bowman-Birk trypsin inhibitor, and Cowpea trypsin inhibitor. The KTIs were first
408 discovered from soybean in 1945 ¹². Since then, a number of trypsin inhibitors have been
409 discovered and characterized from winged bean ^{14,65}, predominantly from seeds, where they
410 are shown to act as insecticidal agents, preventing seed loss during development ⁶⁶. The
411 soybean trypsin inhibitor (STI) gene superfamily has been well studied among *Populus*
412 species, where it has been identified as a rapidly evolving gene family and shown to play
413 multiple roles in anti-herbivory and other stress responses ⁶⁷. Philippe et al ⁶⁴ suggest that the
414 STI gene family has expanded due to repeated gene duplications within poplar relative to
415 other plant species because poplars need strong anti-herbivory actions to maintain their long-
416 lived life cycle. Our study also discovered several *Populus*-specific radiations (subclades A,
417 C, & D; Figure 7).

418

419 In this study, we characterized 28 STI sequences from our CPP34-7 transcripts as well as 20
420 from the Nigerian winged bean transcriptome²³. The vast majority of our STI sequences
421 clustered with *Glycine* in subclade F (Fig. 7), which includes 28 of 32 overall, distinct
422 *Psophocarpus* lineages, 15 of which are corroborated between the CPP34-7 and Chapman
423 transcriptomes. Subclade F includes those proteins originally characterized as KTIs.
424 Expansion of gene family members in *Psophocarpus* can be characterized as lineage
425 (species)-specific or gene-specific (e.g., via tandem gene duplications). The lineage-specific
426 radiations within *Psophocarpus* and *Glycine* may be inflated due to the presence of multiple
427 alleles or alternatively spliced transcripts. Gene-specific amplification of STI family
428 members in poplar is in part due to tandem duplication⁶⁴. KTI genes (with highest sequence
429 similarity to subclade F) are tandemly duplicated within soybean, with at least eight KTI loci
430 linked within 68 kbp on chromosome 8 (between positions 44850000..44918000). Lineage-
431 specific amplification of *Psophocarpus* KTI sequences is evident, and, given the expectation
432 of conserved synteny between soybean and *Psophocarpus*, gene-specific amplification of
433 *Psophocarpus* KTI sequences may be due to recent tandem duplications.

434

435 Besides the classically described KTI genes, several other prominent STI genes are present in
436 our gene tree. Subclade B includes a single contig from our *Psophocarpus* transcriptome,
437 with high sequence similarity to miraculin, a glycoprotein that strongly binds to human taste
438 receptors in the presence of acidic compounds, modifying sour tastes into sweet ones⁶⁸.
439 Miraculin is classified into the STI family and encodes the Kunitz motif but differs from
440 other STI or KTI family members in that it forms a homodimer instead of monomers⁶⁹.
441 Subclade D includes a single *Psophocarpus* contig that has high similarity to alpha-
442 amylase/subtilisin inhibitor proteins known to inhibit the activity of insect α -amylase in

443 *Vigna* species, thus protecting against insect attack ⁷⁰. Subclade E comprises only legume STI
444 gene sequences, including a single paralog from *Psophocarpus* with high sequence similarity
445 to Kunitz-type trypsin inhibitor-like 2 proteins.

446

447 The unequal distribution of *Psophocarpus* STI sequences across the six subclades may be due
448 to tissue specificity, depth of transcriptome reads, amplification of certain gene subfamilies
449 or gene loss over time. As mentioned earlier, most of the winged bean KTIs currently known
450 were characterized from seeds, yet all of these are present in subclade F, in spite of the fact
451 that the CPP34-7 transcriptome did not include seed transcripts, but was sequenced from
452 young leaves, shoots, and buds. This subclade also includes a *Psophocarpus* nodulin
453 (Ptet_Q43325) expressed in nodules of winged bean, likely as a delayed response of the host
454 plant to *Rhizobium* infection ⁷¹. Expression levels of STI genes in winged bean likely differ
455 across plant tissues, as demonstrated in poplar ⁶⁴, and this may be one explanation for
456 unequal distribution of *Psophocarpus* sequences across the gene tree, although inclusion of
457 such tissue-specific genes as nodulins argues against that. Unfortunately, we cannot
458 determine tissue-specific expression of our winged bean STIs due to the fact that our
459 transcriptomes were sequenced from pooled tissue samples. But, given the roles of KTIs in
460 wound and herbivory defense, radiation of KTI genes would be evolutionarily beneficial to
461 large-leaved, highly nutritious plants such as winged bean. Deeper sequencing of the
462 transcriptomes across more tissue types would likely yield other STI gene family members in
463 *Psophocarpus* and provide a more holistic view of STI gene family evolution in the winged
464 bean.

465

466 **Materials and methods**

467

468 *Plant material*

469

470 Seeds of two winged bean (*Psophocarpus tetragonolobus*) genotypes were selected from the
471 United States Department of Agriculture (USDA) Germplasm Resources Information
472 Network (GRIN) seed bank. PI 639033 (CPP37) was field collected in 1999 while PI 491423
473 (CPP34) was donated in 1984, both from Sri Lanka. Seeds were grown to maturity in the
474 greenhouse at Cornell University (Ithaca, NY, USA) for 3 years. Flowering and fruiting were
475 induced by imposing a day length of less than 8 hours. For comparative purposes and to aid
476 in the development of genetic resources for the winged bean, we compared our
477 transcriptomes to an Illumina-based *P. tetragonolobus* transcriptome (SRR1772344) recently
478 published and originally sourced from Nigeria ²³.

479

480 *RNA isolation and library preparation*

481

482 Young leaves, young buds, and young shoots were collected from 3-year old plants into
483 liquid nitrogen to preserve RNA. Total RNA was extracted from each tissue (leaves, shoots,
484 and buds) separately using the Qiagen RNeasy mini kit according to manufacturer
485 instructions. The quality and quantity of each RNA tissue extract was assessed using a 2100
486 Bioanalyzer (Agilent Technologies, Santa Clara, CA, USA). All RNA samples had RIN
487 (RNA integrity number) greater than 9.0 and were used for the analysis. RNA concentration
488 was also quantified using the nanodrop 2000c spectrophotometer (NanoDrop Technologies,
489 Inc., Montchanin, DE, USA). Before cDNA library construction, RNA from tissues for each
490 accession was combined in equal molar amounts so as to allow each tissue equal
491 representation in the final library construct. One microgram (ug) of the pooled tissue total
492 RNA extracts were used for subsequent cDNA library construction of each accession using

493 the Clontech SMARTer cDNA synthesis kit (Clontech Laboratories, Inc., Mountain View,
494 CA, USA) according to manufacturer's instructions but using a 3' SMART CDS Primer IIA
495 modified to 5' --
496 AAGCAGTGGTATCAACGCAGAGTACTTTTTTGTTTTTTCTTTTTTTVN -- 3' which
497 was purchased from IDT (Integrated DNA Technologies, Inc., CA, USA). cDNA libraries
498 were then purified using the PureLink PCR purification kit (Life Technologies, (Invitrogen),
499 Carlsbad, CA, USA) with Buffer HC which removed all fragments less than 300 bp.

500

501 *Transcriptome sequencing*

502

503 Samples were sequenced using single-end 454 pyrosequencing on the Roche 454 Genome
504 Sequencer FLX (Titanium chemistry) at the Brigham Young University Sequencing Center
505 (Provo, UT, USA). Libraries were tagged with multiplex identifier (MID) barcodes to allow
506 multiplexing of four species together over one titer plate. After sequencing, MID adaptors
507 and primers were removed from reads during pre-processing. Preliminary visualization of
508 data was done in FASTQC v. 0.11.3 ⁷².

509

510 *De novo assembly*

511

512 For the CPP34 transcriptome we used the standard flowgram file (SFF) originally generated
513 by the 454 GS FLX sequencer. However, for CPP37, we started with fasta (FNA) and quality
514 value (QUAL) files. We converted the FNA and QUAL files for CPP37 into a single FASTQ
515 file using a python script. We used the FASTX-Toolkit
516 (http://hannonlab.cshl.edu/fastx_toolkit/index.html) to trim and clean the CPP37 reads: we
517 discarded sequences shorter than 50 bp (-l 50) using FASTX CLIPPER and setting of first

518 base of 15 (-f 15) and last base of 800 (-l 800) using FASTA TRIMMER. Finally, the
519 FASTQ Quality Filter was used with minimum quality score of 20 (-q 20) and minimum
520 percent of included bases of 80 (-p 80). In all steps we used the quality score ASCII offset
521 command (-Q33) to denote 454 file format. The quality of output reads after cleaning steps
522 was inspected using FASTQC software v. 0.11.3 ⁷².

523

524 To determine the extent of divergence between our two independently sequenced 454-based
525 transcriptomes, CPP34 and CPP37, we initially assembled each transcriptome independently
526 and explored several contemporary assembly strategies, including Trinity ⁷³, Velvet ⁷⁴, MIRA
527 ⁷⁵, and GS *De Novo* Assembler (aka Newbler, Roche, USA) (see methods and results in
528 Supplementary file 1 online). Our initial findings found fewer than 200 high confidence
529 SNPs between assemblies of CPP34 and CPP37 (SNPs were detected between CPP34 and
530 CPP37 the same way they were assessed between Sri Lankan and Nigerian accessions; see
531 SNP methods below), suggesting a high degree of similarity between these two Sri Lankan
532 accessions. Therefore, for subsequent assemblies and analyses we combined the reads from
533 our two Sri Lankan accessions and produced a single assembly, notated as CPP34-7.

534 Ultimately, we chose to use GS *De Novo* Assembler over the other programs because of the
535 reliable output, comparable contig length, the fact that it considers alternative splicing ⁷⁶, and
536 that it is a program specifically designed for 454 data. Comparisons for several programs in
537 the past showed that it performed best among other de novo assemblers for 454 transcriptome
538 data ⁷⁷. Raw reads from CPP34 and CPP37 were combined by co-assembly within GS *De*
539 *Novo* Assembler v. 2.9 with default settings using a minimum read length of 20, minimum
540 overlap length of 40, minimum overlap identity of 90%, and Isotig threshold of 100.

541

542 *Functional annotation*

543

544 Prior to functional annotation, we identified candidate coding regions and filtered sequences
545 based on a minimum amino acid length of 100 using the TransDecoder program
546 (<https://transdecoder.github.io>) v. 2.0.1 applied to CPP34-7 contigs plus singletons, using the
547 TransDecoder.LongOrfs command. To identify open reading frames (ORFs) with homology
548 to known proteins and to maximize sensitivity for capturing ORFs that may have functional
549 significance, Blastp and Pfam searches were conducted. The Blastp search was done using
550 the Swissprot database with the E-value of 1E-5 and Pfam search was done using HMMER
551 (<http://hmmer.janelia.org>), a biosequence analysis program using profile hidden Markov
552 models and the Pfam database (<http://pfam.xfam.org>). Output files that were generated from
553 the Blastp and Pfam database searches were leveraged by TransDecoder to ensure that
554 peptides with BLAST or domain hits were retained in the set of reported likely coding
555 regions by running the TransDecoder.Predict command. Finally, output of the TransDecoder
556 analysis was used as input for functional annotation using the Blast2GO program⁷⁸. First, we
557 conducted a BLAST search on the output from Transdecoder against the NCBI's
558 nonredundant (nr) database with the E-value of 1E-5 on the Smithsonian Hydra clusters.
559 These BLAST results were then used as input to Blast2GO to assign Gene Ontology (GO)
560 terms to our DNA regions, including biological processes (BP), molecular functions (MF),
561 and cellular components (CC).

562

563 *Sequence similarity with other legumes*

564

565 To compare our complement of genes characterized from our winged bean transcriptome
566 assembly against typical gene assemblies in other legumes, legume species' protein
567 sequences (*Medicago truncatula* Gaertn., *Glycine max* (L.) Merr. (soybean), *Lotus japonicus*

568 (Regel) K.Larsen, *Phaseolus vulgaris* L. (common bean), *Cicer arietinum* L. (chickpea), and
569 *Cajanus cajan* (L.) Millsp. (pigeonpea)) along with *Populus trichocarpa* Torr. & A.Gray ex
570 Hook. and *Arabidopsis thaliana* (L.) Heynh. protein sequences were downloaded from NCBI.
571 BLASTX searches were performed against the CPP34-7 contigs with E-value of 1E-4, and
572 the top hit for each contig was used for further analysis.

573

574 *Transcription factor identification*

575

576 CPP34-7 transcripts were translated to protein sequences for prediction of transcription
577 factors in the assembly. Translated protein sequences were subjected to prediction using
578 PlantTFDB (<http://planttfdb.cbi.pku.edu.cn/>), with further linking the prediction to best hits
579 in Arabidopsis. Since not all transcription factors (TFs) could be predicted in the CPP34-7
580 assembly, we utilized the annotation results of BLASTX searches against legume databases.
581 All identified and predicted transcription factors were further classified into categories.

582

583 *Simple sequence repeat identification*

584

585 To retrieve simple sequence repeat (SSRs; microsatellite) markers and also to design primers,
586 SSR Locator v.1⁷⁹ program was used to detect SSRs across contigs from CPP34-7. A SSR
587 site was defined as a monomer occurring at least 12x with a dimer at least 6x, trimers at least
588 4x, tetra- and pentamers at least 3x, and hexa- to decamers occurring at least 2x. The space
589 between compound SSRs was set to 100 bp and the space between imperfect SSRs to 5 bp.
590 Primers were produced and reported for primary SSRs only.

591

592 *Single nucleotide polymorphism identification*

593

594 To identify SNPs between *Psophocarpus* transcriptomes, we used the transcripts (contigs +
595 singletons) from our combined assembly CPP34-7 as a reference ‘genome’. We extracted
596 singletons from the original reads and concatenated them with the contigs produced by our
597 CPP34-7 assembly. We queried Chapman’s Nigerian, Illumina-based transcriptome²³ against
598 our CPP34-7 reference ‘genome’ using the GUI interface of GS Reference Mapper v. 2.9
599 (454 Life Sciences, Roche, USA) under default settings. We used only high-confidence
600 variants to the reference sequence (454HCDiffs) and further filtered these to those having
601 20x or greater coverage. Lastly, to ensure the highest SNP call quality for use in future
602 research, we followed the method of Schmutz *et al*²¹ and discarded any SNPs where i) the
603 reference or variant involved one or more N’s; and or ii) the reference or variant allele was a
604 single nucleotide insertion or deletion or did not involve a point mutation in the length
605 variant.

606

607 *Kunitz-type trypsin inhibitor gene family analysis*

608

609 To reconstruct a gene tree of the STI superfamily, particularly the KTI gene families, and to
610 understand the evolutionary diversification of this gene superfamily in *Psophocarpus* related
611 to other legumes, we obtained available STI sequences for selected legumes and other
612 angiosperms from the Pfam database (<http://pfam.xfam.org>). In total, 214 accessions were
613 retrieved across *Arabidopsis*, *Populus*, *Medicago*, *Phaseolus*, *Glycine*, and *Psophocarpus*.
614 We downloaded a reference alignment from the Pfam database and used this alignment as a
615 scaffold upon which to align contigs garnered from our transcriptome (see Supplementary
616 file 6 online). We extracted putative *Psophocarpus* STI regions from our transcriptome
617 (CPP34-7) and Chapman’s²³ after blasting against a local BLAST database⁸⁰ based on

618 available STI gene sequences obtained from the Pfam database.

619

620 We combined our extracted contigs with the Pfam STI sequences, converted the open reading
621 frames to amino acid sequences, and aligned in MAFFT v. 7. 245⁸¹. Phylogenetic analysis
622 was conducted in RAxML v. 8.1.24⁸² with 1000 rapid bootstrap inferences and using the
623 best substitution model (LG+G) as determined by Prottest v. 3.4⁸³. Additionally, Bayesian
624 analysis was conducted using MrBayes v. 3.2.6⁸⁴ under the JTT amino acid model
625 "aamodelpr = fixed(jones)" and gamma rates. Two independent Markov Chain Monte Carlo
626 (MCMC) analyses with 12 simultaneous chains and 25 million generations were run for each
627 analysis. Trees were sampled every 10,000 generations and the first 25% of trees were
628 discarded as burn-in. The convergence of MCMC chains was confirmed with Tracer version
629 1.6⁸⁵. All runs and parameters were checked to ensure proper mixing as evidenced by
630 effective sample size (ESS) scores being above 200 and the standard deviation of the split
631 frequencies having dropped below 0.01⁸⁴.

632

633 **Acknowledgements**

634

635 Thanks go to Sue Sherman-Broyles and Jane L. Doyle for help in the greenhouse. Thanks to
636 Matthew Kweskin and Vanessa Gonzalez from the Smithsonian National Museum of Natural
637 History for help with data analysis. Computations were completed on the Smithsonian
638 Institution High Performance Cluster (SI/HPC). This research was supported by grants from
639 the US National Science Foundation to JJD (DEB-0948800) and ANE (DEB-1352217). We
640 acknowledge the support of Sir M Visvesvaraya Institute of Technology (Sir MVIT),
641 Bangalore and Sikkim University, Gangtok for facilities.

642

643 **Authors' contributions**

644

645 All authors contributed to various aspects of this work (ordered by degree of contribution):
646 conceived the study (ANE, NS), aided in study design (ANE, NS, MV, PS), obtained funds
647 for the research (JJD, ANE), coordinated activities (ANE, NS), obtained and grew plants
648 from seed (ANE), extracted RNA and prepared cDNA libraries for 454 sequencing (ANE),
649 conducted bioinformatic analyses (MV, PS, SS, ANE, RC), and contributed to preparation of
650 the manuscript (MV, ANE, NS, PS, JJD, RC). All authors reviewed the manuscript.

651

652 **Availability of Supporting Data**

653

654 Transcriptome datasets supporting the conclusions of this article are available in the NCBI
655 SRA repository under the accession number SRP067662 (raw 454 reads). In addition, several
656 large datasets stemming from analyses of these data are available in Supplementary files
657 online.

658

659 **Competing Financial Interests**

660

661 The authors declare that they have no competing financial interests.

662

663 **References**

664 1 Hymowitz, T. & Boyd, J. Origin, ethnobotany and agricultural potential of the winged
665 bean—*Psophocarpus tetragonolobus*. *Econ. Bot.* **31**, 180-188 (1977).

- 666 2 Klu, G. Induced mutations for accelerated domestication - a case study of winged
667 bean (*Psophocarpus tetragonolobus* (L.) DC). *West African Journal of Applied*
668 *Ecology* **1**, 47-52 (2000).
- 669 3 Harder, D. K. Chromosome counts in *Psophocarpus*. *Kew Bulletin* **47**, 529-534
670 (1992).
- 671 4 Smith, O., Ilori, J. & Onesirosan, P. The proximate composition and nutritive value of
672 the winged bean *Psophocarpus tetragonolobus* (L.) DC for broilers. *Anim. Feed Sci.*
673 *Technol.* **11**, 231-237 (1984).
- 674 5 Amoo, I., Adebayo, O. & Oyeleye, A. Chemical evaluation of winged beans
675 (*Psophocarpus tetragonolobus*), Pitanga cherries (*Eugenia uniflora*) and orchid fruit
676 (Orchid fruit myristica). *African Journal of Food, Agriculture, Nutrition and*
677 *Development* **6**, 1-12 (2006).
- 678 6 Bean, N. R. C.-P. o. t. W. *The winged bean: a high-protein crop for the tropics*.
679 (National Academies, 1975).
- 680 7 Harder, D., Lolema, O. P. M. & Tshisand, M. Uses, nutritional composition, and
681 ecogeography of four species of *Psophocarpus* (Fabaceae, Phaseoleae) in Zaire. *Econ.*
682 *Bot.* **44**, 391-409 (1990).
- 683 8 Ruegg, J. Effects of temperature and water stress on the growth and yield of winged
684 bean (*Psophocarpus tetragonolobus* (L.) DC.). *J. Hortic. Sci.* **56**, 331-338 (1981).
- 685 9 Tan, N. H., Rahim, Z. H., Khor, H. T. & Wong, K. C. Winged bean (*Psophocarpus*
686 *tetragonolobus*) tannin level, phytate content and hemagglutinating activity. *J. Agric.*
687 *Food Chem.* **31**, 916-917 (1983).
- 688 10 Ryan, C. A. Protease inhibitors in plants: genes for improving defenses against insects
689 and pathogens. *Annu. Rev. Phytopathol.* **28**, 425-449 (1990).

- 690 11 Habu, Y., Fukushima, H., Sakata, Y., Abe, H. & Funada, R. A gene encoding a major
691 Kunitz proteinase inhibitor of storage organs of winged bean is also expressed in the
692 phloem of stems. *Plant Mol. Biol.* **32**, 1209-1213 (1996).
- 693 12 Kunitz, M. Crystallization of a trypsin inhibitor from soybean. *Science* **101**, 668-669
694 (1945).
- 695 13 Peyachoknagul, S. *et al.* Sequence and expression of the mRNA encoding the
696 chymotrypsin inhibitor in winged bean (*Psophocarpus tetragonolobus* (L.) DC.).
697 *Plant Mol. Biol.* **12**, 51-58 (1989).
- 698 14 Giri, A. P. *et al.* Identification of potent inhibitors of *Helicoverpa armigera* gut
699 proteinases from winged bean seeds. *Phytochemistry* **63**, 523-532 (2003).
- 700 15 Harding, J., Martin, F. & Kleiman, R. Seed protein and oil yields of the winged bean
701 *Psophocarpus tetragonolobus* in Puerto Rico. *Tropical Agriculture (Trinidad and*
702 *Tobago)* **55**, 307 (1978).
- 703 16 Klu, G., Quaynor-Addy, M., Dinku, E. & Dikumwin, E. in *Joint FAO/IAEA Division*
704 *of Nuclear Techniques in Food and Agriculture, Vienna (Austria)* 15-16
705 (International Atomic Energy Agency, 1989).
- 706 17 Chen, D. *et al.* Genetic diversity evaluation of winged bean (*Psophocarpus*
707 *tetragonolobus* (L.) DC.) using inter-simple sequence repeat (ISSR). *Genet. Resour.*
708 *Crop Evol.* **62**, 823-828 (2015).
- 709 18 Sharma, K. K., Dumbala, S. R. & Bhatnagar-Mathur, P. in *Plant Biotechnol.* 193-
710 207 (Springer, 2014).
- 711 19 Egan, A. N., Schlueter, J. & Spooner, D. M. Applications of next-generation
712 sequencing in plant biology. *Am. J. Bot.* **99**, 175-185 (2012).
- 713 20 Varshney, R. K., Close, T. J., Singh, N. K., Hoisington, D. A. & Cook, D. R. Orphan
714 legume crops enter the genomics era! *Curr. Opin. Plant Biol.* **12**, 202-210 (2009).

- 715 21 Schmutz, J. *et al.* A reference genome for common bean and genome-wide analysis of
716 dual domestications. *Nat. Genet.* **46**, 707-713 (2014).
- 717 22 Lavin, M., Herendeen, P. S. & Wojciechowski, M. F. Evolutionary rates analysis of
718 Leguminosae implicates a rapid diversification of lineages during the Tertiary. *Syst.*
719 *Biol.* **54**, 575-594 (2005).
- 720 23 Chapman, M. A. Transcriptome sequencing and marker development for four
721 underutilized legumes. *Applications in Plant Sciences* **3**, apps. 1400111,
722 doi:doi:10.3732/apps.1400111. (2015).
- 723 24 Nelson, R. J., Naylor, R. L. & Jahn, M. M. The role of genomics research in
724 improvement of “orphan” crops. *Crop Sci.* **44**, 1901-1904 (2004).
- 725 25 Varshney, R. K., Nayak, S. N., May, G. D. & Jackson, S. A. Next-generation
726 sequencing technologies and their implications for crop genetics and breeding. *Trends*
727 *Biotechnol.* **27**, 522-530 (2009).
- 728 26 Graham, I. in *Successful Agricultural Innovation in Emerging Economies: New*
729 *Genetic Technologies for Global Food Production* (eds D.J. Bennet & R.C.
730 Jennings) 95-106 (2013).
- 731 27 Sharpe, A. G. *et al.* Ancient orphan crop joins modern era: gene-based SNP discovery
732 and mapping in lentil. *BMC Genomics* **14**, 192, doi:DOI: 10.1186/1471-2164-14-192
733 (2013).
- 734 28 Hiremath, P. J. *et al.* Large- scale transcriptome analysis in chickpea (*Cicer arietinum*
735 L.), an orphan legume crop of the semi- arid tropics of Asia and Africa. *Plant*
736 *Biotechnol. J.* **9**, 922-931 (2011).
- 737 29 Yang, T. *et al.* Large-scale microsatellite development in grasspea (*Lathyrus sativus*
738 L.), an orphan legume of the arid areas. *BMC Plant Biol.* **14**, 65, doi:DOI:
739 10.1186/1471-2229-14-65 (2014).

- 740 30 Chen, H. *et al.* Transcriptome sequencing of mung bean (*Vigna radiate* L.) genes and
741 the identification of EST-SSR markers. *PloS One* **10**, e0120273, doi:DOI:
742 10.1371/journal.pone.0120273 (2015).
- 743 31 Souframanien, J. & Reddy, K. S. De novo assembly, characterization of immature
744 seed transcriptome and development of genic-SSR markers in black gram [*Vigna*
745 *mungo* (L.) Hepper]. *PloS One* **10**, e0128748 (2015).
- 746 32 FAO, I. *WFP. 2015.* (2015).
- 747 33 Kudapa, H. *et al.* Comprehensive transcriptome assembly of chickpea (*Cicer*
748 *arietinum* L.) using Sanger and next generation sequencing platforms: development
749 and applications. *PLos One* **9**, e86039 (2014).
- 750 34 Sudheesh, S. *et al.* De novo assembly and characterisation of the field pea
751 transcriptome using RNA-Seq. *BMC Genomics* **16**, 611 (2015).
- 752 35 Pop, M. & Salzberg, S. L. Bioinformatics challenges of new sequencing technology.
753 *Trends Genet.* **24**, 142-149 (2008).
- 754 36 Meyer, E. *et al.* Sequencing and de novo analysis of a coral larval transcriptome using
755 454 GSFlx. *BMC Genomics* **10**, 219 (2009).
- 756 37 Libault, M. *et al.* Legume transcription factor genes: what makes legumes so special?
757 *Plant Physiol.* **151**, 991-1001 (2009).
- 758 38 Schmutz, J. *et al.* Genome sequence of the paleopolyploid soybean. *Nature* **463**, 178-
759 183, doi:10.1038/nature08670 (2010).
- 760 39 Wang, Z. *et al.* SoyDB: a knowledge database of soybean transcription factors. *BMC*
761 *Plant Biol.* **10**, 14 (2010).
- 762 40 Guimarães, P. M. *et al.* Global transcriptome analysis of two wild relatives of peanut
763 under drought and fungi infection. *BMC Genomics* **13**, 387 (2012).

- 764 41 Llorca, C. M., Potschin, M. & Zentgraf, U. bZIPs and WRKYs: two large
765 transcription factor families executing two different functional strategies. *Front. Plant*
766 *Sci.* **5**, 10-3389 (2014).
- 767 42 Martín-Trillo, M. & Cubas, P. TCP genes: a family snapshot ten years later. *Trends*
768 *Plant Sci.* **15**, 31-39 (2010).
- 769 43 Ma, J. *et al.* Genome-wide identification and expression analysis of TCP transcription
770 factors in *Gossypium raimondii*. *Scientific Reports* **4**, 6645 (2014).
- 771 44 Yanhui, C. *et al.* The MYB transcription factor superfamily of Arabidopsis:
772 expression analysis and phylogenetic comparison with the rice MYB family. *Plant*
773 *Mol. Biol.* **60**, 107-124 (2006).
- 774 45 Wang, M. L., Barkley, N. A. & Jenkins, T. M. Microsatellite markers in plants and
775 insects. Part I: Applications of biotechnology. *Genes, Genomes and Genomics* **3**, 54-
776 67 (2009).
- 777 46 Varshney, R. K., Graner, A. & Sorrells, M. E. Genic microsatellite markers in plants:
778 features and applications. *Trends Biotechnol.* **23**, 48-55 (2005).
- 779 47 Ellis, J. & Burke, J. EST-SSRs as a resource for population genetic analyses. *Heredity*
780 **99**, 125-132 (2007).
- 781 48 Chankaew, S. *et al.* Detection of genome donor species of neglected tetraploid crop
782 *Vigna reflexo-pilosa* (creole bean), and genetic structure of diploid species based on
783 newly developed EST-SSR markers from azuki bean (*Vigna angularis*). *PLoS One* **9**,
784 e104990, doi:doi:10.1371/journal.pone.0104990 (2014).
- 785 49 Sun, X. *et al.* SSR genetic linkage map construction of pea (*Pisum sativum* L.) based
786 on Chinese native varieties. *The Crop Journal* **2**, 170-174 (2014).
- 787 50 Eujayl, I. *et al.* *Medicago truncatula* EST-SSRs reveal cross-species genetic markers
788 for *Medicago* spp. *Theor. Appl. Genet.* **108**, 414-422 (2004).

- 789 51 Bosamia, T. C., Mishra, G. P., Thankappan, R. & Dobarra, J. R. Novel and stress
790 relevant EST derived SSR markers developed and validated in peanut. *PloS One* **10**,
791 e0129127 (2015).
- 792 52 Mun, J.-H. *et al.* Distribution of microsatellites in the genome of *Medicago*
793 *truncatula*: a resource of genetic markers that integrate genetic and physical maps.
794 *Genetics* **172**, 2541-2555 (2006).
- 795 53 Zhang, L. *et al.* Preference of simple sequence repeats in coding and non-coding
796 regions of *Arabidopsis thaliana*. *Bioinformatics* **20**, 1081-1086 (2004).
- 797 54 Blair, M. W., Torres, M. M., Giraldo, M. C. & Pedraza, F. Development and diversity
798 of Andean-derived, gene-based microsatellites for common bean (*Phaseolus vulgaris*
799 L.). *BMC Plant Biol.* **9**, 100 (2009).
- 800 55 Castro, J. C. *et al.* De novo assembly and functional annotation of *Myrciaria dubia*
801 fruit transcriptome reveals multiple metabolic pathways for L-ascorbic acid
802 biosynthesis. *BMC Genomics* **16**, 997 (2015).
- 803 56 Tóth, G., Gáspári, Z. & Jurka, J. Microsatellites in different eukaryotic genomes:
804 survey and analysis. *Genome Res.* **10**, 967-981 (2000).
- 805 57 Morgante, M., Hanafey, M. & Powell, W. Microsatellites are preferentially associated
806 with nonrepetitive DNA in plant genomes. *Nat. Genet.* **30**, 194-200 (2002).
- 807 58 Rajesh, P. & Muehlbauer, F. J. Discovery and detection of single nucleotide
808 polymorphism (SNP) in coding and genomic sequences in chickpea (*Cicer arietinum*
809 L.). *Euphytica* **162**, 291-300 (2008).
- 810 59 Hyten, D. L. *et al.* High-throughput SNP discovery through deep resequencing of a
811 reduced representation library to anchor and orient scaffolds in the soybean whole
812 genome sequence. *BMC Genomics* **11**, 38 (2010).

- 813 60 Mammadov, J., Aggarwal, R., Buyyarapu, R. & Kumpatla, S. SNP markers and their
814 impact on plant breeding. *International Journal of Plant Genomics* **2012**, Article ID
815 728398, doi:doi:10.1155/2012/728398 (2012).
- 816 61 Coulondre, C., Miller, J. H., Farabaugh, P. J. & Gilbert, W. Molecular basis of base
817 substitution hotspots in *Escherichia coli*. *Nature* **274**, 775-780 (1978).
- 818 62 Agarwal, G. *et al.* Comparative analysis of kabuli chickpea transcriptome with desi
819 and wild chickpea provides a rich resource for development of functional markers.
820 *PLoS One* **7**, e52443 (2012).
- 821 63 Leonforte, A. *et al.* SNP marker discovery, linkage map construction and
822 identification of QTLs for enhanced salinity tolerance in field pea (*Pisum sativum* L.).
823 *BMC Plant Biol.* **13**, 161 (2013).
- 824 64 Philippe, R. N., Ralph, S. G., Külheim, C., Jancsik, S. I. & Bohlmann, J. Poplar
825 defense against insects: genome analysis, full- length cDNA cloning, and
826 transcriptome and protein analysis of the poplar Kunitz- type protease inhibitor
827 family. *New Phytol.* **184**, 865-884 (2009).
- 828 65 Yamamoto, M., Saburo, H. & Ikenaka, T. Amino acid sequences of two trypsin
829 inhibitors from winged bean seeds (*Psophocarpus tetragonolobus* (L) DC.). *J.*
830 *Biochem.* **94**, 849-863 (1983).
- 831 66 Gatehouse, A. M., Hoe, D. S., Flemming, J. E., Hilder, V. A. & Gatehouse, J. A.
832 Biochemical basis of insect resistance in winged bean (*Psophocarpus tetragonolobus*)
833 seeds. *J. Sci. Food Agric.* **55**, 63-74 (1991).
- 834 67 Major, I. T. & Constabel, C. P. Functional analysis of the Kunitz trypsin inhibitor
835 family in poplar reveals biochemical diversity and multiplicity in defense against
836 herbivores. *Plant Physiol.* **146**, 888-903 (2008).

- 837 68 Theerasilp, S. & Kurihara, Y. Complete purification and characterization of the taste-
838 modifying protein, miraculin, from miracle fruit. *J. Biol. Chem.* **263**, 11536-11539
839 (1988).
- 840 69 Takai, A. *et al.* Secretion of miraculin through the function of a signal peptide
841 conserved in the Kunitz-type soybean trypsin inhibitor family. *FEBS Lett.* **587**, 1767-
842 1772 (2013).
- 843 70 Kokiladevi, E., Manickam, A. & Thayumanavan, B. Characterization of alpha-
844 amylase inhibitor in *Vigna sublobata*. *Botanical Bulletin of Academia Sinica* **46**
845 (2005).
- 846 71 Manen, J.-F. *et al.* A nodulin specifically expressed in senescent nodules of winged
847 bean is a protease inhibitor. *Plant Cell* **3**, 259-270 (1991).
- 848 72 Andrews, S. FastQC: A quality control tool for high throughput sequence data.
849 *Reference Source* (2010).
- 850 73 Grabherr, M. G. *et al.* Trinity: reconstructing a full-length transcriptome without a
851 genome from RNA-Seq data. *Nat. Biotechnol.* **29**, 644 (2011).
- 852 74 Zerbino, D. R. & Birney, E. Velvet: algorithms for de novo short read assembly using
853 de Bruijn graphs. *Genome Res.* **18**, 821-829 (2008).
- 854 75 Chevreux, B., Wetter, T. & Suhai, S. in *German conference on bioinformatics.* 45-
855 56.
- 856 76 Mundry, M., Bornberg-Bauer, E., Sammeth, M. & Feulner, P. G. Evaluating
857 characteristics of de novo assembly software on 454 transcriptome data: a simulation
858 approach. *PloS One* **7**, e31410 (2012).
- 859 77 Kumar, S. & Blaxter, M. L. Comparing de novo assemblers for 454 transcriptome
860 data. *BMC Genomics* **11**, 571, doi:DOI: 10.1186/1471-2164-11-571 (2010).

- 861 78 Conesa, A. *et al.* Blast2GO: a universal tool for annotation, visualization and analysis
862 in functional genomics research. *Bioinformatics* **21**, 3674-3676,
863 doi:10.1093/bioinformatics/bti610 (2005).
- 864 79 Da Maia, L. C. *et al.* SSR Locator: Tool for simple sequence repeat discovery
865 integrated with primer design and PCR simulation. *International Journal of Plant*
866 *Genomics* (2008).
- 867 80 McGinnis, S. & Madden, T. L. BLAST: at the core of a powerful and diverse set of
868 sequence analysis tools. *Nucleic Acids Res.* **32**, W20-W25 (2004).
- 869 81 Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version
870 7: improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772-780 (2013).
- 871 82 Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis
872 of large phylogenies. *Bioinformatics* **30**, 1312-1313 (2014).
- 873 83 Darriba, D., Taboada, G. L., Doallo, R. & Posada, D. ProtTest 3: fast selection of
874 best-fit models of protein evolution. *Bioinformatics* **27**, 1164-1165 (2011).
- 875 84 Ronquist, F. *et al.* MrBayes 3.2: efficient Bayesian phylogenetic inference and model
876 choice across a large model space. *Syst. Biol.* **61**, 539-542 (2012).
- 877 85 Rambaut, A., Suchard, M. A., Xie, D. & Drummond, A. J. Tracer v1.5, available from
878 <http://beast.bio.ed.ac.uk/Tracer>. (2013).

879

880 **Figure 1: Summary of Gene Annotation analysis.** Zero Hit refers to those in BLAST step
881 without hits.

882

883 **Figure 2: Gene Ontology classifications of Winged bean annotated transcripts.** Numbers
884 indicate the number of sequences associated with the particular GO term in each category.

885

886 **Figure 3: % Identity of CPP34-7 contigs against legume protein databases.**

887

888 **Figure 4: Legume sequence similarity analysis.** Relative numbers of contigs that had
889 significant sequence similarity by species for CPP34-7 contigs.

890

891 **Figure 5: Transcription factor family analysis.** Number of transcription factors determined
892 within the CPP34-7 assembly by transcription factor family.

893

894 **Figure 6: Results of microsatellite SSR analyses.** (A) Distribution of the 2,994 perfect
895 SSRs across different repeat size classes. Distribution of the number of repeats for (B) dimers
896 (C) trimers (D) tetramers (E) pentamers and (F) hexamers.

897

898 **Figure 7: Gene Tree of Kunitz Trypsin Inhibitor Gene Family.** Non-legume sequences:
899 *Arabidopsis thaliana* (At; black), *Populus trichocarpa* (Ptri; black). Legume sequences:
900 *Medicago truncatula* (Mt; pink), *Phaseolus vulgaris* (Pv; orange), *Glycine max* (Gm; green),
901 and *Psophocarpus tetragonolobus* from Pfam database (Ptet; navy blue), Chapman (2015)
902 transcriptome (Pt_c; royal blue), and CPP34-7 (Pt_isotig/contig; aqua blue). Sequence
903 notation is species abbreviation followed by Pfam accession number, contig/isotig number, or
904 read, followed by the range of amino acids used in the alignment. Numbers at the nodes are
905 posterior probability values and bootstrap supports. Subclades A-F are labeled. *Psophocarpus*
906 clades are indicated by arrows or blue banding. Tree rooted arbitrarily at an *Arabidopsis*
907 clade.

908

909 **Table 1.** Sequencing and assembly metrics for independent and combined assemblies using
 910 *GS De Novo Assembler*.

Accessions	Genotype CPP34	Genotype CPP37	Combined Assembly (CPP34-7)
Number of raw reads	371,271	433,486	804,757
Number of bases (bp)	191,598,691	213,386,165	404,984,856
Number of reads post-filtering	369,820 (99.6%)	334,639 (77.2%)	704,459 (87.53%)
Number of bases post-filtering	136,943,216 (71.47%)	92,126,948 (43.17%)	178,911,104 (44.17%)
Number of reads aligned	277,351 (50.42%)	259,324 (63.04%)	435,897 (61.88%)
Number of contigs / bp	10,675 / 6,142,297	8,465 / 5,070,585	16,115 / 13,552,130
Avg. contig size (bp)	837	823	875
N50 (bp)	836	842	889
Longest contig (bp)	4,902	3,014	4,667
Number of singletons / bp	62,602 / 22,081,798	63,795 / 23,540,672	81,126 / 28,663,213
Number of transcripts / bp (contigs + singletons)	73,277 / 28,224,095	72,260 / 28,611,257	97,241 / 42,215,343

911

912

913 **Table 2.** Distribution of di- and trinucleotide repeat motif types in winged bean and
 914 comparison with *Arabidopsis*.

Dinucleotide Repeat Composition	Number of transcripts	Percentage of Winged bean di- Repeats	Winged bean Rank	Percentage of Arabidopsis di- Repeats
AC/CA/GT/TG	22	8.6	3	8
AG/GA/CT/TC	199	77.7	1	83
AT/TA	35	13.7	2	8.8
CG/GC	0	0	4	0.14
Total	256	100		100

Trinucleotide Repeat Composition	Number of transcripts	Percentage of Winged bean tri- Repeats	Winged bean Rank	Arabidopsis Rank
AAC/ACA/CAA/GTT/TGT/TTG	134	11.4	4	3
AAG/AGA/GAA/CTT/TCT/TTC	343	29.1	1	1
AAT/ATA/TAA/TTA/TAT/ATT	58	4.9	8	5
ACC/CAC/CCA/GGT/GTG/TGG	118	10.0	5	4
ACG/CGA/GAC/CGT/GTC/TCG	35	3.0	9	9
ACT/CTA/TAC/AGT/TAG/GTA	13	1.1	10	8
AGC/CAG/GCA/TGC/CTG/GCT	136	11.5	3	7
AGG/GGA/GAG/TCC/CTC/CCT	118	10.0	6	6
ATC/CAT/TCA/GAT/ATG/TGA	164	13.9	2	2
CCG/CGC/GCC/GGC/GCG/CGG	59	5.0	7	10
Total	1,178	100		

915

916 **Table 3.** Results of single nucleotide polymorphism (SNP) detection between Sri Lankan and
 917 Nigerian genotypes by degree of confidence.

918

Reads	95%	96%	97%	98%	99%	100%	Total
Contigs	43	68	76	94	50	2,552	2,883
Singletons	74	88	93	52	28	1,972	2,307
Total	117	156	169	146	78	4,524	5,190

919

920