

# SCIENTIFIC REPORTS



OPEN

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

Received: 04 April 2016

Accepted: 14 June 2016

Published: 30 June 2016

Mohammad Vatanparast<sup>1</sup>, Prateek Shetty<sup>2</sup>, Ratan Chopra<sup>3</sup>, Jeff J. Doyle<sup>4</sup>, N. Sathyanarayana<sup>5</sup> & Ashley N. Egan<sup>1</sup>

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

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

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

<sup>1</sup>US National Herbarium (US), Department of Botany, Smithsonian Institution-NMNH, 10th and Constitution Ave, Washington DC, 20013, USA. <sup>2</sup>Department of Plant Biology, Michigan State University, 612 Wilson Road, Room 166, East Lansing, MI, 48824, USA. <sup>3</sup>United States Department of Agriculture, Agriculture Research Service, 3810 4th St., Lubbock, TX, 79415, USA. <sup>4</sup>Section of Plant Breeding & Genetics, School of Integrative Plant Science, Cornell University, 412 Mann Library, Ithaca, NY, 14853, USA. <sup>5</sup>Department of Botany, Sikkim University, 5th Mile, Tadong, Gangtok, Sikkim, 737102, India. Correspondence and requests for materials should be addressed to N.S. (email: nsathyanarayana@cus.ac.in) or A.N.E. (email: egana@si.edu)

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

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

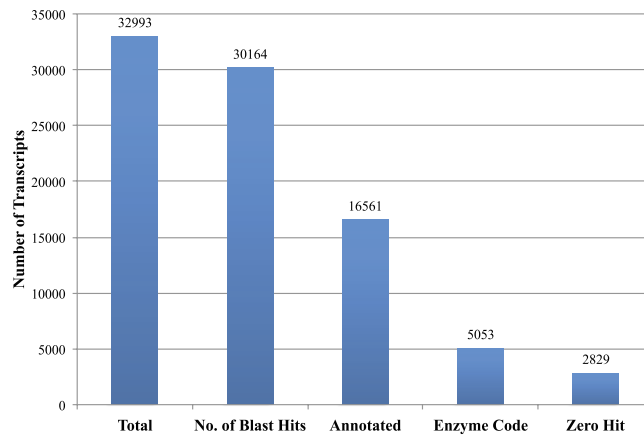
Since the 1975 publication by the US National Academy of Sciences of *The Winged Bean: A High Protein Crop for the Tropics*<sup>6</sup>, considerable effort has been focused on studying the nutritional quality and climatic and ecological tolerances of the plant<sup>7,8</sup>. Winged bean reportedly possesses anti-nutritional factors such as phytolectins, cyanogenic glycosides, tannins, lectins, flatulence factors, and saponins<sup>9</sup>. However, processing using moist heat or soaking has been shown to safely eliminate these substances. Research efforts concerning such anti-nutritional components have yielded significant knowledge concerning trypsin, a serine protease that acts to hydrolyze proteins as part of vertebrate digestion, and trypsin inhibitors, proteins that stop the action of trypsin, thereby interfering with digestion. It has been suggested that trypsin inhibitors play a role in protecting plant tissues against the action of bacterial proteases at the colonization site of pathogenic bacteria<sup>10</sup>. In addition, studies show involvement in defense against insects that suck the phloem sap and against bacteria that invade upon wounding<sup>11</sup>. In biomedical research, these modes of action have made trypsin and trypsin inhibitors vital components of molecular cell research where they are widely used in cell culture to detach cells from tissue culture plates. Since their first discovery in soybeans in 1945<sup>12</sup>, other Kunitz-type trypsin inhibitors have been discovered and characterized from winged bean<sup>13,14</sup>, predominantly from seeds.

It is hard to find another high rainfall-adapted tropical legume with as many desirable characteristics as winged bean<sup>1</sup>. However, much needs to be done in terms of breeding efforts, especially to develop self-supporting, determinate cultivars bearing large numbers of relatively small pods having nutritious seeds and tubers, and cultivars resistant to biotic and abiotic stresses. Considerable variability for growth vigor and quantitative characters such as protein and oil content as well as photoperiodic responses has been recorded<sup>15</sup>. Several beneficial mutants were recovered during the 1970s and 80's through mutation breeding experiments<sup>16</sup>. However, a recent study using inter-simple sequence repeat (ISSR) markers reported low genetic diversity among the winged bean germplasm collected from different parts of the world<sup>17</sup>. With the advent of genomic tools such as molecular markers, genetic maps, etc., the genetic improvement of underutilized crops has been greatly facilitated, enabling the development of improved genotypes or varieties with enhanced trait values<sup>18</sup>. In the case of winged bean, studies on genomic resource development for enabling basic and applied research on genetics, evolution, ecology and molecular breeding programs are lacking, yet the advent of genomic technologies provides significant prospects for improvement<sup>19</sup>. Transcriptome sequencing is cost-effective and a valuable method for efficient and rapid identification of molecular markers in resource poor plant species<sup>20</sup>.

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

## Results

**Sequencing and *De novo* assembly of winged bean transcriptomes.** Pyrosequencing of two Sri Lankan accessions produced comparable sequence output, where genotype CPP34 produced a total of 369,820 single-end reads comprising 136,943,216 bp with an average read length of 574 bp and genotype CPP37 produced a total of 334,639 single-end reads comprising 92,126,948 bp with an average read length of 565 bp (Table 1). Using read count as a proxy, the depth of sequencing across our contigs was similar for the independent *de novo* assemblies, ranging from one to 4,953 reads, with an average read depth of 25 reads per contig for CPP34 and ranging from one to 3,972 reads with an average read depth of 30 reads per contig for CPP37. Comparison of transcripts from the CPP34 and CPP37 independent assemblies (Supplementary file 1, inclusive of Tables S1–S3 and Fig. S1) found fewer than 200 high-confidence SNPs between them (data not shown), equating to



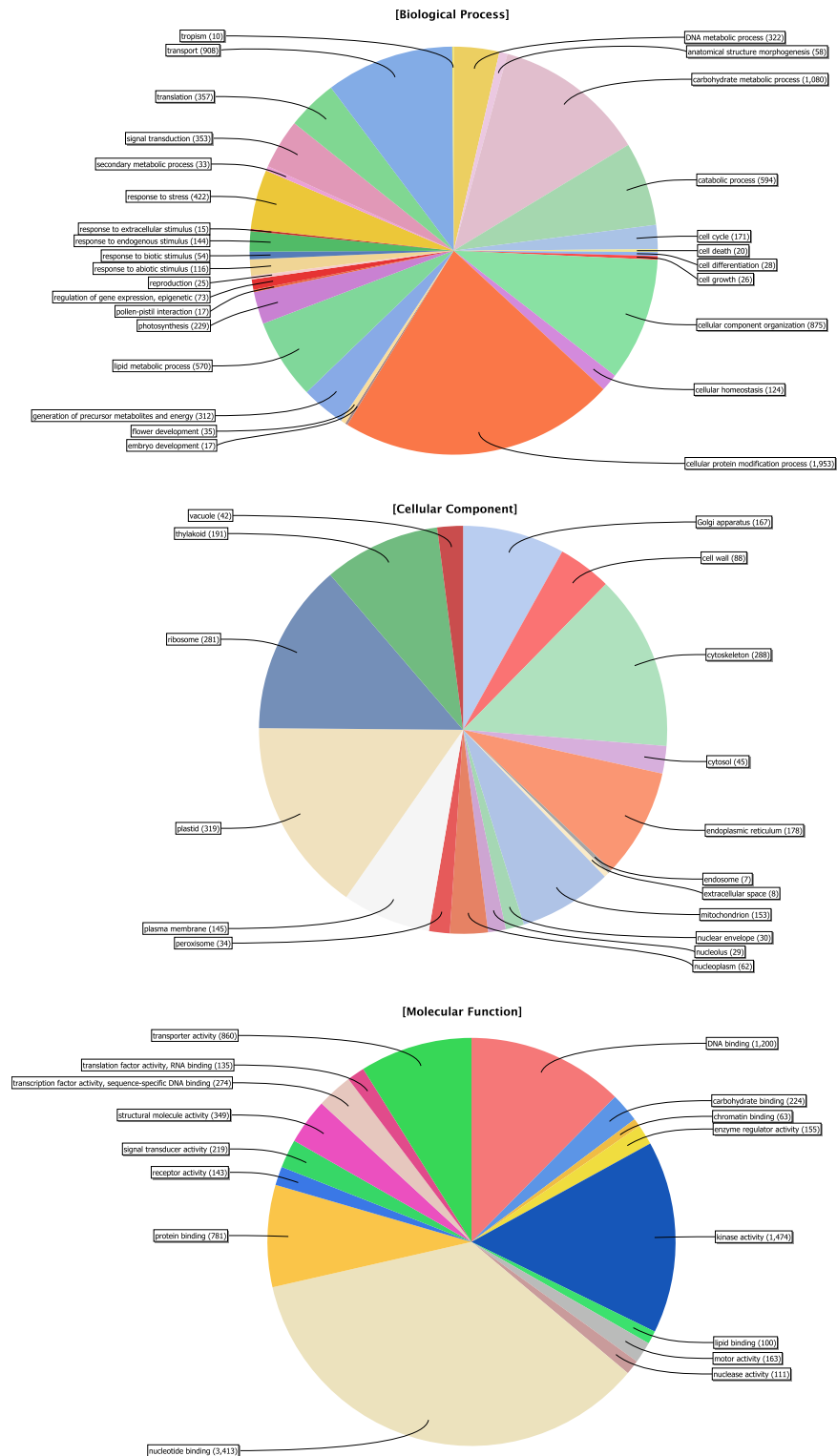
**Figure 1. Summary of gene annotation analysis.** Zero Hit refers to those in BLAST step without hits.

approximately one SNP every 150,000 bp. Therefore, reads from the independently sequenced accessions were combined and co-assembled. For the combined assembly (CPP34-7), this translated to 704,459 reads comprising 229,070,164 bases from both accessions (Table 1). Because 454 pyrosequencing produces comparatively long reads (300–800 bp long), unassembled reads, here notated as singletons post-assembly, may potentially represent full-length mRNA transcripts. In order to not lose potential information, singletons of the CPP34-7 were extracted and appended to the final assembly of CPP34-7 and used in the Gene Ontology (GO) and SNP analyses.

**Functional annotation & legume sequence similarity.** For the GO analysis, the combined assembly of CPP34-7 was used with inclusion of singletons (16,115 contigs plus 81,126 singletons, Table 1). Using a total of 97,241 transcripts, TransDecoder could track 33,038 transcripts against BLAST and Pfam databases. Of these 33,038 transcripts, BLAST searches against NCBI's nr database retrieved 32,993 transcripts with hits (see Supplementary file 2), discarding 45 transcripts that had zero hits in NCBI. Therefore, 64,248 (66%) of our original 97,241 transcripts did not hit any known gene or DNA region in NCBI and Pfam databases, of which 62,783 were singletons. Thus, 79% of singletons were discarded in the BLAST searching steps due to a lack of annotation. Of the 32,993 transcripts with BLAST hits, the GO analysis determined GO ID and enzyme code (EC) assignments for 16,561 (50.1%) with full or partial annotations (Fig. 1 in text, and see Supplementary file 2). Of the 16,561 annotated transcripts, 5,053 have predicted functions (EC codes). Overall, 2,829 transcripts were not functionally annotated by Blast2GO (zero hits) of which 1,932 (68%) corresponded to singletons. Participation of genes in a particular biological process and molecular function are shown in Fig. 2. Several transcripts were assigned to more than one GO term; therefore, the total number of GO terms obtained for our dataset was higher than the total number of transcripts. In total, 47,178 GO terms were retrieved, with 46.2%, 37% and 16.8%, corresponding to the molecular functions (MF), biological processes (BP), and cellular components (CC) categories, respectively. In the MF category, nucleotide binding (number of sequences = 3,413), kinase activity (1,474) and DNA binding (1,200) had the highest number of assigned sequences. In the BP category, cellular protein modification (1,953), carbohydrate metabolic processes (1,080) and transport (908) were the majority and in the case of CC, genes involved in the plastid (319), cytoskeleton (288) and ribosome (281) activities were highly represented (Fig. 2).

A comparison of our assembled contigs against other legume NCBI protein sequence databases from chickpea (*Cicer arietinum* L.), pigeon pea (*Cajanus cajan* (L.) Huth), soybean (*Glycine max* (L.) Merr.), common bean (*Phaseolus vulgaris* L.), *Medicago truncatula* Gaertn., and *Lotus japonicus* (Regel) K.Larsen using the BLASTX program from NCBI showed that 15,558 of 16,115 (96.5%) contigs from the CPP34-7 assembly had significant sequence similarity to sequences in one or more legume protein databases. About 90.5% of the 16,115 contigs had  $\geq 80\%$  sequence identity (Fig. 3). The majority of the contigs (57.3%) were most similar to *G. max* (Fig. 4), a finding that, at first glance, seems to contradict that expected based on evolutionary relationships of legume lineages, but is likely due to the relative over-representation of genes within the soybean genome due to i) recent whole genome duplication and ii) a much higher level and standard of annotation and gene discovery relative to other legume genomes. Differences in evolutionary rates across lineages may also impact this outcome. In relation to *Phaseolus vulgaris*, it is known that *Phaseolus* has a higher mutation rate than *Glycine* and related lineages<sup>21,22</sup>, which could increase the divergence, and thereby decrease the best-BLAST hits, of *Psophocarpus* against *Phaseolus* relative to *Glycine*. However, this explanation is invoked with caution given that it assumes similar relative rates between *Glycine* and *Psophocarpus*, information that is beyond the scope of this project.

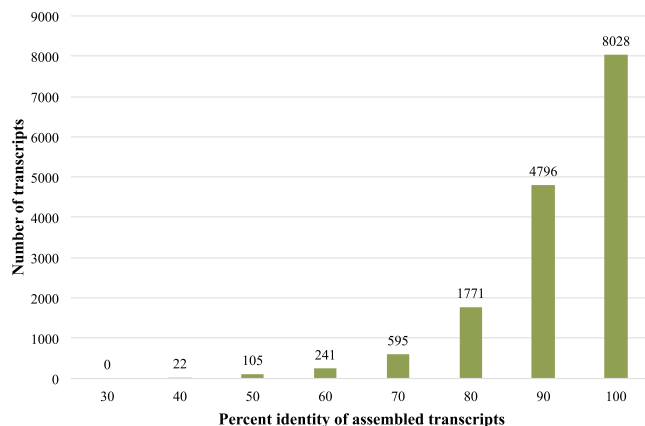
**Identification of transcription factors.** In the overall GO analysis, 274 transcripts were annotated as transcription factors (TFs) (Fig. 2). Of the 16,115 contigs, 176 putative winged bean transcription factor genes, distributed in at least ten families, were identified representing 1.1% of winged bean transcripts, which were assigned to different categories. Among these, basic leucine zipper (bZIP; 32), Teosinte-Branched1/Cycloidea/PCF (TCP; 19), MADS (17), MYB (11) and WRKY (9) were among the top five categories (Fig. 5.). The overall distribution of transcription factor encoding transcripts among the various known protein families is very similar



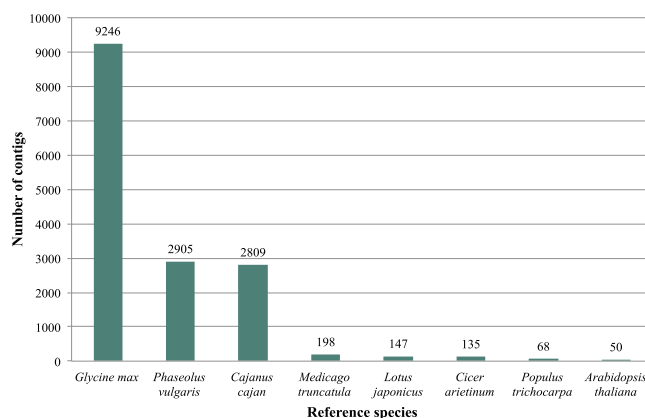
**Figure 2. Gene ontology classifications of winged bean annotated transcripts.** Numbers indicate the number of sequences associated with the particular GO term in each category.

to that of soybean and other legumes. However, almost all families showed minor species specific differences (for example, bZIP, MYB, WRKY etc.) with regard to TF gene families reported for *Lotus*, *Medicago* and *Glycine max*.

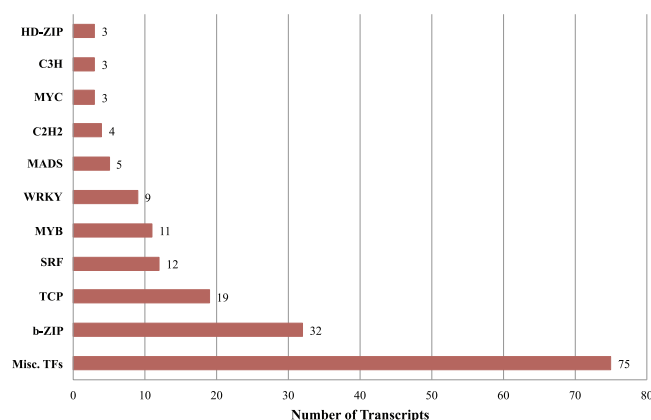
**Identification of simple sequence repeats.** The SSR analysis detected 10,984 perfect SSRs, 13 imperfect SSRs, and 1,959 compound SSRs, for a total of 12,956 SSRs (see Supplementary file 3). Of the 10,984 primary SSRs, 57 were adenine (A: 30) or thymine (T: 27) monomers with at least 13 repeats. These were assumed to



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

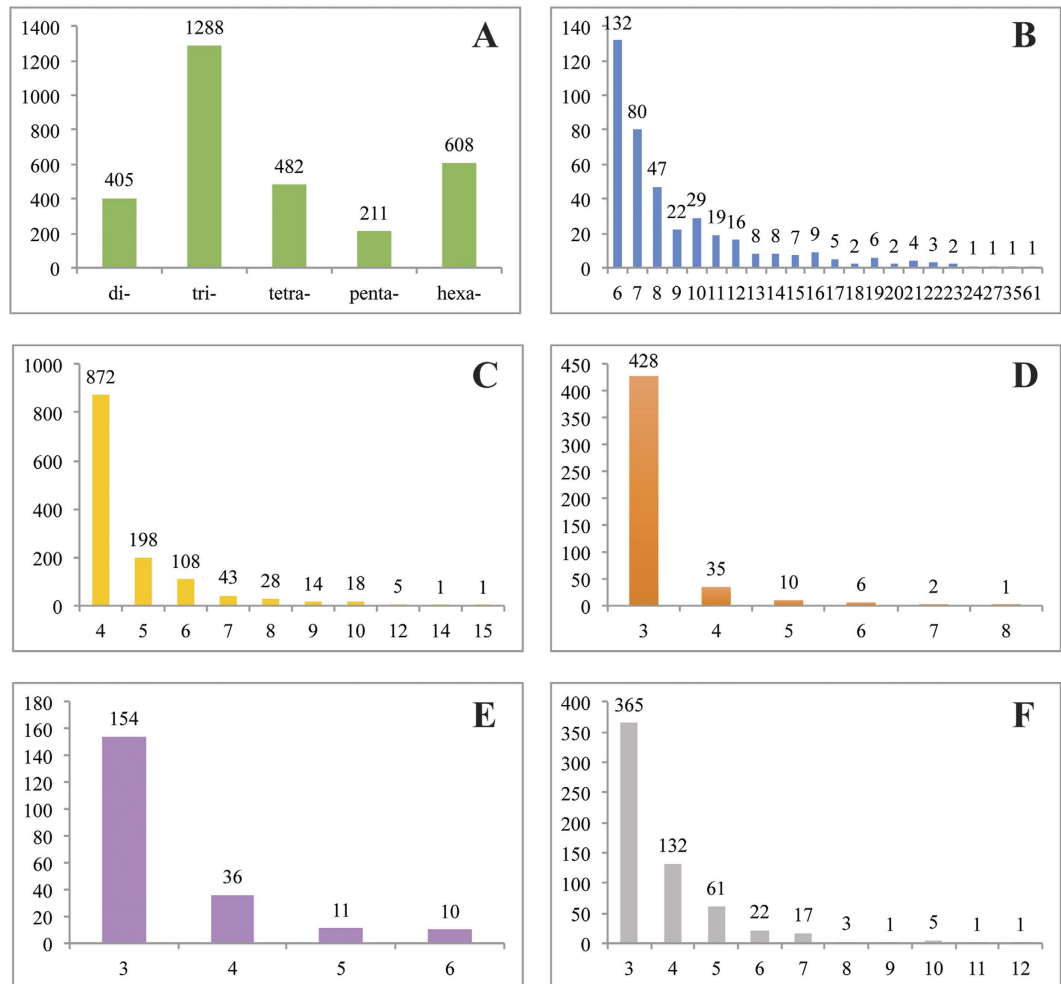


**Figure 4.** Legume sequence similarity analysis. Relative numbers of contigs that had significant sequence similarity by species for CPP34-7 contigs.



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

represent remnants of mRNA poly-A tails and were thus removed prior to primer prediction. No runs of 12 or more cytosine or guanine monomeric repeats were found. Nearly three-quarters of the remaining 10,927 perfect SSRs (7,933) were hexamers with only two repeats. Although these 12-mers may be useful as linkage markers, the low number of repeat units would likely take these out of the microsatellite category. The remaining 2,994 perfect SSRs were distributed across di-, tri-, tetra-, penta-, and hexamer SSRs (Fig. 6) and were used for primer creation. The majority (63%) of SSRs were detected in the tri- and hexamer categories (Fig. 6A). In general, the number of SSRs detected in each size category decreased with increasing repeat number (Fig. 6B–F). Primers



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

were successfully created for 2,594 SSRs with product sizes ranging from 100 to 280 bp (see Supplementary file 4). Analysis of the primed SSRs showed bias towards certain di- and tri- repeat type motifs (Table 2).

**Single nucleotide polymorphism discovery.** GS Reference Mapper mapped 87.7% of reads from Chapman<sup>23</sup> onto the CPP34-7 reference ‘genome’ which consisted of the 97,241 transcripts. Of the 14,571,393 bp of mapped reads, we identified 113,757 SNPs with >95% confidence from the 454HCdiffs file (available upon request), suggesting a SNP frequency of one in 128 bp of coding regions. Interestingly, the majority of high-confidence SNPs were found within singletons (91,686; 80.6%) vs. contigs (22,071; 19.4%), a higher percentage than expected given that singletons make up 67.9% of total transcript length. As a conservative measure, we filtered SNPs based on allele frequency from >95–100% confidence levels and those having >20× coverage (Table 3), producing a total of 13,091 SNPs distributed across 10,176 transcripts of which 5,196 (39.7%) were from contigs, representing 1 SNP every 1,113 bp. The subsequent increase in the proportion of SNPs within contigs is expected in this case given that more highly expressed genes will be more likely to be represented by >20× coverage and are most likely to assemble into contigs. Lastly, we removed all single nucleotide indels (7,665 of the 13,091) and those length variants that involved insertions or deletions of one or more nucleotides alone (i.e. those without point mutations involved in the length variants), resulting in a high-confidence set of 5,190 SNPs, 96% of which are one-to-one point mutations (see Supplemental file 5). Within the 5,190 SNPs, 151 unique SNP patterns were found and 211 (4%) SNPs were length variants involving one or more point mutation within the length variant. Of the 4,979 one-to-one polymorphisms, 3,433 (68.9%) were transitions and 1,546 (31.1%) were transversions, producing a transition:transversion ratio of 2.22.

**Kunitz-type trypsin inhibitor gene family analysis.** We identified 28 contigs from CPP34-7 and 20 contigs from the Chapman<sup>23</sup> transcriptome assembly corresponding to the Kunitz trypsin inhibitor (KTI) gene family within the *Psophocarpus* transcriptome (see Supplementary file 6). Due to the large number of paralogues in each species, there is no obvious criterion available for rooting this tree, so it was rooted with the largest clade

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		

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

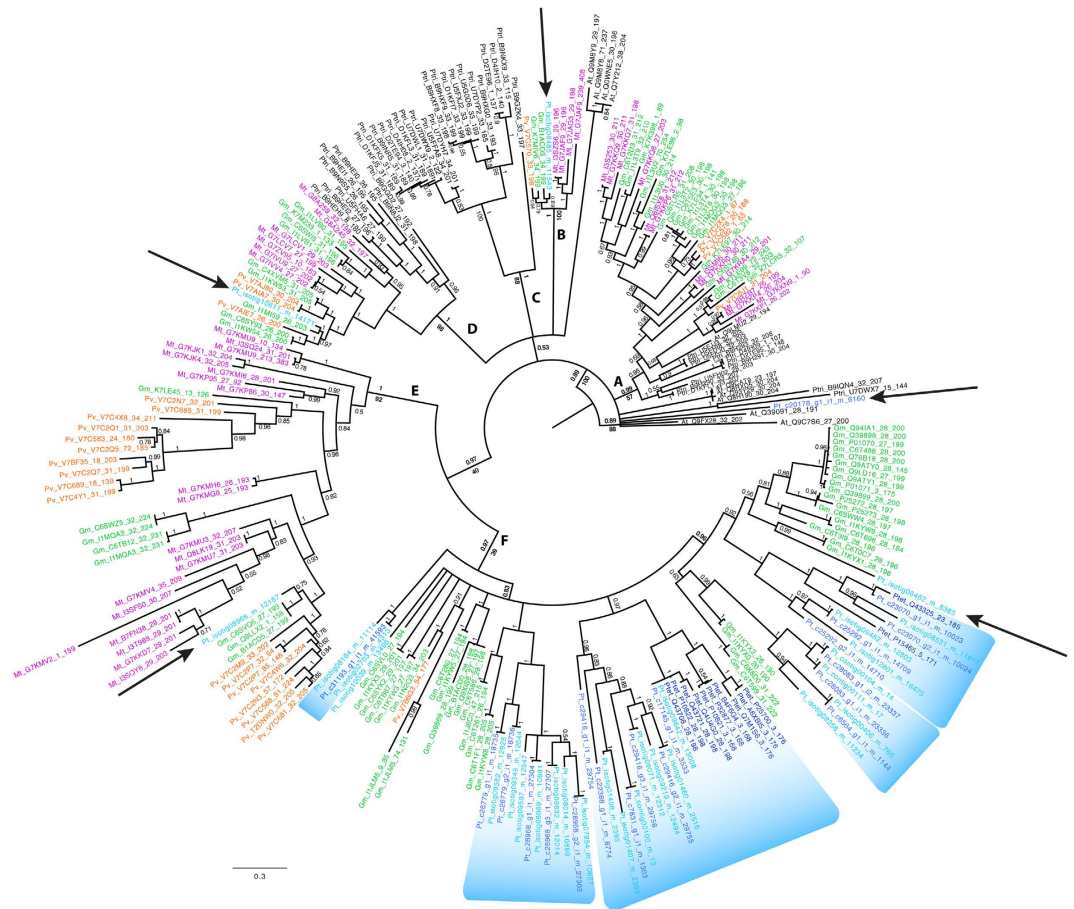
Reads	95%	96%	97%	98%	99%	100%	Total
Contigs	43	68	76	94	50	2,552	2,883
Singletons	74	88	93	52	28	1,972	2,307
Total	117	156	169	146	78	4,524	5,190

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

of non-legume sequences, a clade of *Arabidopsis* sequences. Given this rooting, the Bayesian soybean trypsin inhibitor (STI) gene tree has a polytomous backbone and suggests six distinct subclades based on relatively high posterior probability and bootstrap support, here labeled as A-F (Fig. 7 in text; and see Supplementary file 7). The dominant feature of the tree (regardless of rooting) is a lack of clear orthologous relationships across taxa, with evidence of lineage-specific amplification of STI and KTI genes in each species. For example, subclade A comprises two clades made up of only *Populus trichocarpa* Torr. & A. Gray sequences and an *Arabidopsis thaliana* (L.) Heynh. STI member as well as a number of clades containing *Glycine*, *Phaseolus*, or *Medicago* gene family members, but with no *Psophocarpus* sequences included, whereas subclade C comprises *Populus* sequences only, illustrating a major intra-specific STI radiation (Fig. 7). The vast majority of *Psophocarpus* sequences cluster in clade F, along with many *Glycine* and a single *Phaseolus* sequence. Of the *Psophocarpus* sequences in subclade F, 15 contigs are paired between CPP34-7 and Chapman, forming sister groups that likely represent the same gene in each transcriptome, whereas 13 are unique (Fig. 7). Subclade F illustrates lineage-specific KTI expansion in both *Psophocarpus* and *Glycine*. All *Psophocarpus* sequences obtained from the Pfam or NCBI databases were nested within subclade F, where the majority of the Pfam sequences appeared as monophyletic clades with a contig each from CPP34-7 and Chapman nested therein (Fig. 7).

## Discussion

The legumes represent the third largest family of the flowering plants, many of which are important sources of food, fodder, oil, fiber and medicines. However, with the exception of common pulse crops such as soybean, common bean, etc., a large number of legumes have remained underutilized due to poorly developed infrastructure, especially for genetic and genomic resources<sup>24</sup>. The advent of genomic technologies has brightened the prospects for such orphan crops<sup>20,25,26</sup>, with recent research focusing on lentil (*Lens culinaris* Medik.<sup>27</sup>), chickpea<sup>28</sup>, grass pea (*Lathyrus sativus* L.<sup>29</sup>), and a number of *Vigna* species<sup>30,31</sup>, among others. Winged bean represents a promising alternative to protein-rich soybean for tropical regions of the world that house nearly 40% of the world population, of which nearly one third is protein deficient, and many of whom are women and children<sup>32</sup>. Genomics assisted breeding and enabling biotechnologies that stem from it offer significant promise for targeted genetic improvement of nutritional and other quality traits in winged bean, thus aiding in the development of a low input, high quality legume-based protein diet for these parts of the world. Our combined assembly presents a genetic resource that can be mined for future genetic improvement and plant breeding initiatives. This paper reports development of genetic resources, including molecular markers, in winged bean, in addition to insights into the



**Figure 7. Gene tree of Kunitz trypsin inhibitor gene family.** Non-legume sequences: *Arabidopsis thaliana* (At; black), *Populus trichocarpa* (Ptri; black). Legume sequences: *Medicago truncatula* (Mt; pink), *Phaseolus vulgaris* (Pv; orange), *Glycine max* (Gm; green), and *Psophocarpus tetragonolobus* from Pfam database (Ptet; navy blue), Chapman (2015) transcriptome (Pt\_c; royal blue), and CPP34-7 (Pt\_isotig/contig; aqua blue). Sequence notation is species abbreviation followed by Pfam accession number, contig/isotig number, or read, followed by the range of amino acids used in the alignment. Numbers at the nodes are posterior probability values and bootstrap supports. Subclades A–F are labeled. *Psophocarpus* clades are indicated by arrows or blue banding. Tree rooted arbitrarily at an *Arabidopsis* clade.

divergence of the Kunitz-type trypsin inhibitors, which are important anti-nutritive agents in winged bean and other legumes.

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

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

Because transcription factors play important roles in regulating plant functions, we paid particular attention to their number and distribution within winged bean and in relation to other legumes. Several TF gene families



are preserved across different plant genera, indicating conserved gene regulatory machinery in plants, as has been shown in legumes previously<sup>37</sup>. In this study, we found that 2.4% of the total transcripts are putative transcription factors according to GO analyses, a percentage much lower than the estimated 12% found in soybean (based on ~46,430 protein-coding genes)<sup>38</sup>. However, Libault *et al.*<sup>37</sup> estimated the number of TF-encoding genes across a number of species and found soybean to have 3–4× the number of TFs relative to *Medicago truncatula* or *Lotus japonicus*, likely due to recent whole genome duplication in soybean. If we compare the estimated number of TFs in *M. truncatula* (1,473)<sup>37</sup> against the number of putative protein-coding genes in *M. truncatula* (~66,000; phytozome v11; <https://phytozome.jgi.doe.gov/pz/portal.html>), we come out with a much more similar estimate (2.2%). However, this is likely an underestimate and a shifting target as the annotation of *M. truncatula* is ongoing.

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

Simple sequence repeats, or microsatellite markers, have long been used for genetic diversity analyses and plant breeding efforts, largely due to their highly polymorphic, co-dominant nature, prevalence throughout the genome, ease of use, and cost-effectiveness<sup>45</sup>. Because they originate in coding regions, SSRs derived from genes have increased amplification success in related species, are useful for assessing functional diversity and for marker-assisted selection, and can act as anchor markers for evolutionary and comparative mapping studies<sup>46,47</sup>. While some research has suggested that SSRs derived from coding regions are less polymorphic than their anonymous counterparts<sup>47</sup>, numerous population genetic, evolutionary, and plant breeding studies have found them to have adequate, if not higher, levels of polymorphism within legumes<sup>48,49</sup>.

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

Certain repeat motif types were more prevalent than others in our set of primed SSRs (Table 2), a not-uncommon finding that has been documented in legumes previously<sup>52</sup>. Zhang *et al.*<sup>53</sup> first documented the bias of microsatellites to AG and AAG motifs in *Arabidopsis*, also noting differences of SSR distributions between 5' and 3' untranslated and coding regions, and correlation between trinucleotide repeat motifs and codon usage. In winged bean, the SSR repeat motif type (AG/GA/TC/CT)<sub>n</sub> represented the majority (77.7%) of all dinucleotide repeats, while motif types (AT/TA)<sub>n</sub> and (AC/CA/GT/TG)<sub>n</sub> comprised 13.7% and 8.6%, respectively. Our distribution and ranking of dinucleotide repeat motifs mirrors that in *Arabidopsis* (Table 2). The bias towards the repeat type containing AG and against that of GC has also been found in other plants, including *Phaseolus*<sup>54</sup>, *Myrciaria dubia* (Kunth) McVaugh<sup>55</sup> and across eukaryotes<sup>56</sup>. Past research has suggested that AG motifs are most prevalent in 5' untranslated regions<sup>52,57</sup> and possibly are involved in transcription and regulation<sup>53</sup>. As mentioned earlier, 25.7% of primed genic-SSR transcripts are unannotated, some of which may correspond to 5' untranslated regions where AG motif types are more prevalent. The frequency of trinucleotide repeat motif types was biased towards AAG in our set of primed SSRs, with this motif type comprising 29.1% of the 10 trinucleotide types, followed by the ATC motif type, comprising 13.9%. The ranking of these and other motifs closely resembles that of *Arabidopsis* (Table 2), with the first two most prevalent motifs the same<sup>53</sup>.

SNPs provide another means of assessing genetic variation and, although less polymorphic than SSRs, are abundant and easily obtained via high-throughput sequencing. For example, Rajesh and Muehlbauer<sup>58</sup> estimated SNP frequency to be one in 66 bp in coding regions and one in 71 bp in genomic regions of chickpea<sup>58</sup>. In another study, Hyten *et al.*<sup>59</sup> reported 7,000–25,000 predicted SNPs through deep resequencing of soybean by a whole genome sequence approach. In this study, we discovered more than 5,190 high-confidence SNPs between our Sri Lankan samples and the geographically separated Nigerian genotype<sup>23</sup>. SNP markers identified in this study can be used in quantitative trait loci (QTL) mapping, generating linkage maps, genotyping and breeding studies. Validation of SNPs determined herein is beyond the scope of this paper, nevertheless, this list presents a significant resource for future work in plant breeding and genetic diversity assessment<sup>60</sup> and marks the first SNP markers discovered to date in *Psophocarpus*.

Our high-confidence SNP set included 4,979 one-to-one SNPs (those without length variants and involving changes between a single nucleotide position), equating to a transition:transversion (ts:tv) ratio of 2.22. This bias is commonly observed across SNPs throughout a genome, resulting from rampant methyl cytosine to uracil mutations<sup>61</sup>. Similar ratios were found across SNPs in other legumes<sup>62,63</sup>. In total, 44% of SNPs identified were found in singletons, a proportion not unexpected given that 68% of transcript read length is in singletons. But the very fact that alignable and putatively homologous singletons were found across the geographically separated and independently sequenced genotypes provides vindication for their inclusion in transcriptome characterizations,

at least for 454 data. However, caution is warranted due to the ‘singleness’ of the unassembled read acting as a reference sequence.

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

In this study, we characterized 28 STI sequences from our CPP34-7 transcripts as well as 20 from the Nigerian winged bean transcriptome<sup>23</sup>. The majority of our STI sequences clustered with *Glycine* in subclade F (Fig. 7), which includes 28 of 32 overall, distinct *Psophocarpus* lineages, 15 of which are corroborated between the CPP34-7 and Chapman transcriptomes. Subclade F includes those proteins originally characterized as KTIs. Expansion of gene family members in *Psophocarpus* can be characterized as lineage (species)-specific or gene-specific (e.g., via tandem gene duplications). The lineage-specific radiations within *Psophocarpus* and *Glycine* may be inflated due to the presence of multiple alleles or alternatively spliced transcripts. Gene-specific amplification of STI family members in poplar is in part due to tandem duplication<sup>64</sup>. KTI genes (with highest sequence similarity to subclade F) are tandemly duplicated within soybean, with at least eight KTI loci linked within 68 kbp on chromosome 8 (between positions 44850000..44918000). Lineage-specific amplification of *Psophocarpus* KTI sequences is evident, and, given the expectation of conserved synteny between soybean and *Psophocarpus*, gene-specific amplification of *Psophocarpus* KTI sequences may be due to recent tandem duplications.

Besides the classically described KTI genes, several other prominent STI genes are present in our gene tree. Subclade B includes a single contig from our *Psophocarpus* transcriptome, with high sequence similarity to miraculin, a glycoprotein that strongly binds to human taste receptors in the presence of acidic compounds, modifying sour tastes into sweet ones<sup>68</sup>. Miraculin is classified into the STI family and encodes the Kunitz motif but differs from other STI or KTI family members in that it forms a homodimer instead of monomers<sup>69</sup>. Subclade D includes a single *Psophocarpus* contig that has high similarity to alpha-amylase/subtilisin inhibitor proteins known to inhibit the activity of insect  $\alpha$ -amylase in *Vigna* species, thus protecting against insect attack<sup>70</sup>. Subclade E comprises only legume STI gene sequences, including a single paralog from *Psophocarpus* with high sequence similarity to Kunitz-type trypsin inhibitor-like 2 proteins.

The unequal distribution of *Psophocarpus* STI sequences across the six subclades may be due to tissue specificity, depth of transcriptome reads, amplification of certain gene subfamilies or gene loss over time. As mentioned earlier, most of the winged bean KTIs currently known were characterized from seeds, yet all of these are present in subclade F, in spite of the fact that the CPP34-7 transcriptome did not include seed transcripts, but was sequenced from young leaves, shoots, and buds. This subclade also includes a *Psophocarpus* nodulin (Ptet\_Q43325) expressed in nodules of winged bean, likely as a delayed response of the host plant to *Rhizobium* infection<sup>71</sup>. Expression levels of STI genes in winged bean likely differ across plant tissues, as demonstrated in poplar<sup>64</sup>, and this may be one explanation for unequal distribution of *Psophocarpus* sequences across the gene tree, although inclusion of such tissue-specific genes as nodulins argues against that. Unfortunately, we cannot determine tissue-specific expression of our winged bean STIs due to the fact that our transcriptomes were sequenced from pooled tissue samples. But, given the roles of KTIs in wound and herbivory defense, radiation of KTI genes would be evolutionarily beneficial to large-leaved, highly nutritious plants such as winged bean. Deeper sequencing of the transcriptomes across more tissue types would likely yield other STI gene family members in *Psophocarpus* and provide a more holistic view of STI gene family evolution in the winged bean.

## Materials and Methods

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

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

microgram ( $\mu\text{g}$ ) of the pooled tissue total RNA extracts were used for subsequent cDNA library construction of each accession using the Clontech SMARTer cDNA synthesis kit (Clontech Laboratories, Inc., Mountain View, CA, USA) according to manufacturer's instructions but using a 3' SMART CDS Primer IIA modified to 5'-AAGCAGTGGTATCAACGCAGAGTACTTTTTTGTCTTTTTTCTTTTTTNN-3' which was purchased from IDT (Integrated DNA Technologies, Inc., CA, USA). cDNA libraries were then purified using the PureLink PCR purification kit (Life Technologies, (Invitrogen), Carlsbad, CA, USA) with Buffer HC which removed all fragments less than 300 bp.

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

**De novo assembly.** For the CPP34 transcriptome we used the standard flowgram file (SFF) originally generated by the 454 GS FLX sequencer. However, for CPP37, we started with fasta (FNA) and quality value (QUAL) files. We converted the FNA and QUAL files for CPP37 into a single FASTQ file using a python script. We used the FASTX-Toolkit ([http://hannonlab.cshl.edu/fastx\\_toolkit/index.html](http://hannonlab.cshl.edu/fastx_toolkit/index.html)) to trim and clean the CPP37 reads: we discarded sequences shorter than 50 bp (-l 50) using FASTX CLIPPER and setting of first base of 15 (-f 15) and last base of 800 (-l 800) using FASTA TRIMMER. Finally, the FASTQ Quality Filter was used with minimum quality score of 20 (-q 20) and minimum percent of included bases of 80 (-p 80). In all steps we used the quality score ASCII offset command (-Q 33) to denote 454 file format. The quality of output reads after cleaning steps was inspected using FASTQC software v. 0.11.3<sup>72</sup>.

To determine the extent of divergence between our two independently sequenced 454-based transcriptomes, CPP34 and CPP37, we initially assembled each transcriptome independently and explored several contemporary assembly strategies, including Trinity<sup>73</sup>, Velvet<sup>74</sup>, MIRA<sup>75</sup>, and GS *De Novo* Assembler (aka Newbler, Roche, USA) (see methods and results in Supplementary file 1). Our initial findings found fewer than 200 high confidence SNPs between assemblies of CPP34 and CPP37 (SNPs were detected between CPP34 and CPP37 the same way they were assessed between Sri Lankan and Nigerian accessions; see SNP methods below), suggesting a high degree of similarity between these two Sri Lankan accessions. Therefore, for subsequent assemblies and analyses we combined the reads from our two Sri Lankan accessions and produced a single assembly, notated as CPP34-7. Ultimately, we chose to use GS *De Novo* Assembler over the other programs because of the reliable output, comparable contig length, the fact that it considers alternative splicing<sup>76</sup>, and that it is a program specifically designed for 454 data. Comparisons for several programs in the past showed that it performed best among other *de novo* assemblers for 454 transcriptome data<sup>77</sup>. Raw reads from CPP34 and CPP37 were combined by co-assembly within GS *De Novo* Assembler v. 2.9 with default settings using a minimum read length of 20, minimum overlap length of 40, minimum overlap identity of 90%, and Isotig threshold of 100.

**Functional annotation.** Prior to functional annotation, we identified candidate coding regions and filtered sequences based on a minimum amino acid length of 100 using the TransDecoder program (<https://transdecoder.github.io>) v. 2.0.1 applied to CPP34-7 contigs plus singletons, using the TransDecoder.LongOrfs command. To identify open reading frames (ORFs) with homology to known proteins and to maximize sensitivity for capturing ORFs that may have functional significance, Blastp and Pfam searches were conducted. The Blastp search was done using the Swissprot database with the E-value of 1E-5 and Pfam search was done using HMMER (<http://hmmer.janelia.org>), a biosequence analysis program using profile hidden Markov models and the Pfam database (<http://pfam.xfam.org>). Output files that were generated from the Blastp and Pfam database searches were leveraged by TransDecoder to ensure that peptides with BLAST or domain hits were retained in the set of reported likely coding regions by running the TransDecoder.Predict command. Finally, output of the TransDecoder analysis was used as input for functional annotation using the Blast2GO program<sup>78</sup>. First, we conducted a BLAST search on the output from Transdecoder against the NCBI's nonredundant (nr) database with the E-value of 1E-5 on the Smithsonian Hydra clusters. These BLAST results were then used as input to Blast2GO to assign Gene Ontology (GO) terms to our DNA regions.

**Sequence similarity with other legumes.** To compare our complement of genes characterized from our winged bean transcriptome assembly against typical gene assemblies in other legumes, legume species' protein sequences (*Medicago truncatula*, *Glycine max*, *Lotus japonicus*, *Phaseolus vulgaris*, *Cicer arietinum*, and *Cajanus cajan*) along with *Populus trichocarpa* and *Arabidopsis thaliana* protein sequences were downloaded from NCBI. BLASTX searches were performed against the CPP34-7 contigs with E-value of 1E-4, and the top hit for each contig was used for further analysis.

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

**Simple sequence repeat identification.** To retrieve simple sequence repeat (SSRs; microsatellite) markers and also to design primers, SSR Locator v.1<sup>79</sup> program was used to detect SSRs across contigs from CPP34-7. A

SSR site was defined as a monomer occurring at least 12× with a dimer at least 6×, trimers at least 4×, tetra- and pentamers at least 3×, and hexa- to decamers occurring at least 2×. The space between compound SSRs was set to 100 bp and the space between imperfect SSRs to 5 bp. Primers were produced and reported for primary SSRs only.

**Single nucleotide polymorphism identification.** To identify SNPs between *Psophocarpus* transcriptomes, we used the transcripts (contigs + singletons) from our combined assembly CPP34-7 as a reference ‘genome’. We extracted singletons from the original reads and concatenated them with the contigs produced by our CPP34-7 assembly. We queried Chapman’s Nigerian, Illumina-based transcriptome<sup>23</sup> against our CPP34-7 reference ‘genome’ using the GUI interface of GS Reference Mapper v. 2.9 (454 Life Sciences, Roche, USA) under default settings. We used only high-confidence variants to the reference sequence (454HCDiffs) and further filtered these to those having 20× or greater coverage. Lastly, to ensure the highest SNP call quality for use in future research, we followed the method of Schmutz *et al.*<sup>21</sup> and discarded any SNPs where i) the reference or variant involved one or more N’s; and or ii) the reference or variant allele was a single nucleotide insertion or deletion or did not involve a point mutation in the length variant.

**Kunitz-type trypsin inhibitor gene family analysis.** To reconstruct a gene tree of the STI superfamily, particularly the KTI gene families, and to understand the evolutionary diversification of this gene superfamily in *Psophocarpus* related to other legumes, we obtained available STI sequences for selected legumes and other angiosperms from the Pfam database (<http://pfam.xfam.org>). In total, 214 accessions were retrieved across *Arabidopsis*, *Populus*, *Medicago*, *Phaseolus*, *Glycine*, and *Psophocarpus*. We downloaded a reference alignment from the Pfam database and used this alignment as a scaffold upon which to align contigs garnered from our transcriptome (see Supplementary file 6). We extracted putative *Psophocarpus* STI regions from our transcriptome (CPP34-7) and Chapman’s<sup>23</sup> after blasting against a local BLAST database<sup>80</sup> based on available STI gene sequences obtained from the Pfam database.

We combined our extracted contigs with the Pfam STI sequences, converted the open reading frames to amino acid sequences, and aligned in MAFFT v. 7. 245<sup>81</sup>. Phylogenetic analysis was conducted in RAXML v. 8.1.24<sup>82</sup> with 1000 rapid bootstrap inferences and using the best substitution model (LG + G) as determined by Prottest v. 3.4<sup>83</sup>. Additionally, Bayesian analysis was conducted using MrBayes v. 3.2.6<sup>84</sup> under the JTT amino acid model “aamodelpr = fixed(jones)” and gamma rates. Two independent Markov Chain Monte Carlo (MCMC) analyses with 12 simultaneous chains and 25 million generations were run for each analysis. Trees were sampled every 10,000 generations and the first 25% of trees were discarded as burn-in. The convergence of MCMC chains was confirmed with Tracer version 1.6<sup>85</sup>. All runs and parameters were checked to ensure proper mixing as evidenced by effective sample size (ESS) scores being above 200 and the standard deviation of the split frequencies having dropped below 0.01<sup>84</sup>.

## References

- Hymowitz, T. & Boyd, J. Origin, ethnobotany and agricultural potential of the winged bean—*Psophocarpus tetragonolobus*. *Econ. Bot.* **31**, 180–188 (1977).
- Klu, G. Induced mutations for accelerated domestication - a case study of winged bean (*Psophocarpus tetragonolobus* (L.) DC). *West African Journal of Applied Ecology* **1**, 47–52 (2000).
- Harder, D. K. Chromosome counts in *Psophocarpus*. *Kew Bull.* **47**, 529–534 (1992).
- Smith, O., Ilori, J. & Onesiroso, P. The proximate composition and nutritive value of the winged bean *Psophocarpus tetragonolobus* (L.) DC for broilers. *Anim. Feed Sci. Technol.* **11**, 231–237 (1984).
- Amoo, I., Adebayo, O. & Oyeleye, A. Chemical evaluation of winged beans. (*Psophocarpus tetragonolobus*), Pitanga cherries (*Eugenia uniflora*) and orchid fruit (Orchid fruit myristica). *Afr. J. Food Agric. Nutr. Dev.* **6**, 1–12 (2006).
- Bean, N. R. C.-P. o. t. W. *The winged bean: a high-protein crop for the tropics*. (National Academies, 1975).
- Harder, D., Lolema, O. P. M. & Tshisand, M. Uses, nutritional composition, and ecogeography of four species of *Psophocarpus* (Fabaceae, Phaseoleae) in Zaire. *Econ. Bot.* **44**, 391–409 (1990).
- Ruegg, J. Effects of temperature and water stress on the growth and yield of winged bean (*Psophocarpus tetragonolobus* (L.) DC). *J. Horticult. Sci.* **56**, 331–338 (1981).
- Tan, N. H., Rahim, Z. H., Khor, H. T. & Wong, K. C. Winged bean (*Psophocarpus tetragonolobus*) tannin level, phytate content and hemagglutinating activity. *J. Agric. Food Chem.* **31**, 916–917 (1983).
- Ryan, C. A. Protease inhibitors in plants: genes for improving defenses against insects and pathogens. *Annu. Rev. Phytopathol.* **28**, 425–449 (1990).
- Habu, Y., Fukushima, H., Sakata, Y., Abe, H. & Funada, R. A gene encoding a major Kunitz proteinase inhibitor of storage organs of winged bean is also expressed in the phloem of stems. *Plant Mol. Biol.* **32**, 1209–1213 (1996).
- Kunitz, M. Crystallization of a trypsin inhibitor from soybean. *Science* **101**, 668–669 (1945).
- Peyachoknagul, S. *et al.* Sequence and expression of the mRNA encoding the chymotrypsin inhibitor in winged bean (*Psophocarpus tetragonolobus* (L.) DC). *Plant Mol. Biol.* **12**, 51–58 (1989).
- Giri, A. P. *et al.* Identification of potent inhibitors of *Helicoverpa armigera* gut proteinases from winged bean seeds. *Phytochemistry* **63**, 523–532 (2003).
- Harding, J., Martin, F. & Kleiman, R. Seed protein and oil yields of the winged bean *Psophocarpus tetragonolobus* in Puerto Rico. *Tropical Agriculture (Trinidad and Tobago)* **55**, 307 (1978).
- Klu, G., Quaynor-Addy, M., Dinku, E. & Dikumwin, E. In *Joint FAO/IAEA Division of Nuclear Techniques in Food and Agriculture, Vienna (Austria)* 15–16 (International Atomic Energy Agency, 1989).
- Chen, D. *et al.* Genetic diversity evaluation of winged bean (*Psophocarpus tetragonolobus* (L.) DC.) using inter-simple sequence repeat (ISSR). *Genet. Resour. Crop Evol.* **62**, 823–828 (2015).
- Sharma, K. K., Dumbala, S. R. & Bhatnagar-Mathur, P. In *Plant Biotechnol.* 193–207 (Springer, 2014).
- Egan, A. N., Schlueter, J. & Spooner, D. M. Applications of next-generation sequencing in plant biology. *Am. J. Bot.* **99**, 175–185 (2012).
- Varshney, R. K., Close, T. J., Singh, N. K., Hoisington, D. A. & Cook, D. R. Orphan legume crops enter the genomics era! *Curr. Opin. Plant Biol.* **12**, 202–210 (2009).
- Schmutz, J. *et al.* A reference genome for common bean and genome-wide analysis of dual domestications. *Nat. Genet.* **46**, 707–713 (2014).

22. Lavin, M., Herendeen, P. S. & Wojciechowski, M. F. Evolutionary rates analysis of Leguminosae implicates a rapid diversification of lineages during the Tertiary. *Syst. Biol.* **54**, 575–594 (2005).
23. Chapman, M. A. Transcriptome sequencing and marker development for four underutilized legumes. *Appl. Plant Sci.* **3**, apps.1400111, doi: 10.3732/apps.1400111. (2015).
24. Nelson, R. J., Naylor, R. L. & Jahn, M. M. The role of genomics research in improvement of “orphan” crops. *Crop Sci.* **44**, 1901–1904 (2004).
25. Varshney, R. K., Nayak, S. N., May, G. D. & Jackson, S. A. Next-generation sequencing technologies and their implications for crop genetics and breeding. *Trends Biotechnol.* **27**, 522–530 (2009).
26. Graham, I. In *Successful Agricultural Innovation in Emerging Economies: New Genetic Technologies for Global Food Production* (eds Bennet, D. J. & Jennings, R. C.) 95–106 (2013).
27. Sharpe, A. G. *et al.* Ancient orphan crop joins modern era: gene-based SNP discovery and mapping in lentil. *BMC Genomics* **14**, 192, doi: 10.1186/1471-2164-14-192 (2013).
28. Hiremath, P. J. *et al.* Large-scale transcriptome analysis in chickpea (*Cicer arietinum* L.), an orphan legume crop of the semi-arid tropics of Asia and Africa. *Plant Biotechnol. J.* **9**, 922–931 (2011).
29. Yang, T. *et al.* Large-scale microsatellite development in grasspea (*Lathyrus sativus* L.), an orphan legume of the arid areas. *BMC Plant Biol.* **14**, 65, doi: 10.1186/1471-2229-14-65 (2014).
30. Chen, H. *et al.* Transcriptome sequencing of mung bean (*Vigna radiate* L.) genes and the identification of EST-SSR markers. *PLoS One* **10**, e0120273, doi: 10.1371/journal.pone.0120273 (2015).
31. Soufmanien, J. & Reddy, K. S. *De novo* assembly, characterization of immature seed transcriptome and development of genic-SSR markers in black gram [*Vigna mungo* (L.) Hepper]. *PLoS One* **10**, e0128748 (2015).
32. FAO, I. *WFP 2015*. (2015).
33. Kudapa, H. *et al.* Comprehensive transcriptome assembly of chickpea (*Cicer arietinum* L.) using Sanger and next generation sequencing platforms: development and applications. *PLoS One* **9**, e86039 (2014).
34. Sudheesh, S. *et al.* *De novo* assembly and characterisation of the field pea transcriptome using RNA-Seq. *BMC Genomics* **16**, 611 (2015).
35. Pop, M. & Salzberg, S. L. Bioinformatics challenges of new sequencing technology. *Trends Genet.* **24**, 142–149 (2008).
36. Meyer, E. *et al.* Sequencing and *de novo* analysis of a coral larval transcriptome using 454 GSFLx. *BMC Genomics* **10**, 219 (2009).
37. Libault, M. *et al.* Legume transcription factor genes: what makes legumes so special? *Plant Physiol.* **151**, 991–1001 (2009).
38. Schmutz, J. *et al.* Genome sequence of the paleopolyploid soybean. *Nature* **463**, 178–183, doi: 10.1038/nature08670 (2010).
39. Wang, Z. *et al.* SoyDB: a knowledge database of soybean transcription factors. *BMC Plant Biol.* **10**, 14 (2010).
40. Guimarães, P. M. *et al.* Global transcriptome analysis of two wild relatives of peanut under drought and fungi infection. *BMC Genomics* **13**, 387 (2012).
41. Llorca, C. M., Potschin, M. & Zentgraf, U. bZIPs and WRKYs: two large transcription factor families executing two different functional strategies. *Front. Plant Sci.* **5**, 10–3389 (2014).
42. Martín-Trillo, M. & Cubas, P. TCP genes: a family snapshot ten years later. *Trends Plant Sci.* **15**, 31–39 (2010).
43. Ma, J. *et al.* Genome-wide identification and expression analysis of TCP transcription factors in *Gossypium raimondii*. *Sci. Rep.* **4**, 6645 (2014).
44. Yanhui, C. *et al.* The MYB transcription factor superfamily of Arabidopsis: expression analysis and phylogenetic comparison with the rice MYB family. *Plant Mol. Biol.* **60**, 107–124 (2006).
45. Wang, M. L., Barkley, N. A. & Jenkins, T. M. Microsatellite markers in plants and insects. Part I: Applications of biotechnology. *G3* **3**, 54–67 (2009).
46. Varshney, R. K., Graner, A. & Sorrells, M. E. Genic microsatellite markers in plants: features and applications. *Trends Biotechnol.* **23**, 48–55 (2005).
47. Ellis, J. & Burke, J. EST-SSRs as a resource for population genetic analyses. *Heredity* **99**, 125–132 (2007).
48. Chankaew, S. *et al.* Detection of genome donor species of neglected tetraploid crop *Vigna reflexo-pilosa* (creole bean), and genetic structure of diploid species based on newly developed EST-SSR markers from azuki bean (*Vigna angularis*). *PLoS One* **9**, e104990, doi: 10.1371/journal.pone.0104990 (2014).
49. Sun, X. *et al.* SSR genetic linkage map construction of pea (*Pisum sativum* L.) based on Chinese native varieties. *Crop J.* **2**, 170–174 (2014).
50. Eujayl, I. *et al.* *Medicago truncatula* EST-SSRs reveal cross-species genetic markers for *Medicago* spp. *Theor. Appl. Genet.* **108**, 414–422 (2004).
51. Bosamia, T. C., Mishra, G. P., Thankappan, R. & Dobarra, J. R. Novel and stress relevant EST derived SSR markers developed and validated in peanut. *PLoS One* **10**, e0129127 (2015).
52. Mun, J.-H. *et al.* Distribution of microsatellites in the genome of *Medicago truncatula*: a resource of genetic markers that integrate genetic and physical maps. *Genetics* **172**, 2541–2555 (2006).
53. Zhang, L. *et al.* Preference of simple sequence repeats in coding and non-coding regions of *Arabidopsis thaliana*. *Bioinformatics* **20**, 1081–1086 (2004).
54. Blair, M. W., Torres, M. M., Giraldo, M. C. & Pedraza, F. Development and diversity of Andean-derived, gene-based microsatellites for common bean (*Phaseolus vulgaris* L.). *BMC Plant Biol.* **9**, 100 (2009).
55. Castro, J. C. *et al.* *De novo* assembly and functional annotation of *Myrciaria dubia* fruit transcriptome reveals multiple metabolic pathways for L-ascorbic acid biosynthesis. *BMC Genomics* **16**, 997 (2015).
56. Tóth, G., Gáspári, Z. & Jurka, J. Microsatellites in different eukaryotic genomes: survey and analysis. *Genome Res.* **10**, 967–981 (2000).
57. Morgante, M., Hanafey, M. & Powell, W. Microsatellites are preferentially associated with nonrepetitive DNA in plant genomes. *Nat. Genet.* **30**, 194–200 (2002).
58. Rajesh, P. & Muehlbauer, F. J. Discovery and detection of single nucleotide polymorphism (SNP) in coding and genomic sequences in chickpea (*Cicer arietinum* L.). *Euphytica* **162**, 291–300 (2008).
59. Hyten, D. L. *et al.* High-throughput SNP discovery through deep resequencing of a reduced representation library to anchor and orient scaffolds in the soybean whole genome sequence. *BMC Genomics* **11**, 38 (2010).
60. Mammadov, J., Aggarwal, R., Buyyarapu, R. & Kumpatla, S. SNP markers and their impact on plant breeding. *Int. J. Plant Genomics* **2012**, Article ID 728398, doi: 10.1155/2012/728398 (2012).
61. Coulondre, C., Miller, J. H., Farabaugh, P. J. & Gilbert, W. Molecular basis of base substitution hotspots in *Escherichia coli*. *Nature* **274**, 775–780 (1978).
62. Agarwal, G. *et al.* Comparative analysis of kabuli chickpea transcriptome with desi and wild chickpea provides a rich resource for development of functional markers. *PLoS One* **7**, e52443 (2012).
63. Leonforte, A. *et al.* SNP marker discovery, linkage map construction and identification of QTLs for enhanced salinity tolerance in field pea (*Pisum sativum* L.). *BMC Plant Biol.* **13**, 161 (2013).
64. Philippe, R. N., Ralph, S. G., Külheim, C., Jancsik, S. I. & Bohlmann, J. Poplar defense against insects: genome analysis, full-length cDNA cloning, and transcriptome and protein analysis of the poplar Kunitz-type protease inhibitor family. *New Phytol.* **184**, 865–884 (2009).

65. Yamamoto, M., Saburo, H. & Ikenaka, T. Amino acid sequences of two trypsin inhibitors from winged bean seeds (*Psophocarpus tetragonolobus* (L) DC). *J. Biochem.* **94**, 849–863 (1983).
66. Gatehouse, A. M., Hoe, D. S., Flemming, J. E., Hilder, V. A. & Gatehouse, J. A. Biochemical basis of insect resistance in winged bean (*Psophocarpus tetragonolobus*) seeds. *J. Sci. Food Agric.* **55**, 63–74 (1991).
67. Major, I. T. & Constabel, C. P. Functional analysis of the Kunitz trypsin inhibitor family in poplar reveals biochemical diversity and multiplicity in defense against herbivores. *Plant Physiol.* **146**, 888–903 (2008).
68. Theerasilp, S. & Kurihara, Y. Complete purification and characterization of the taste-modifying protein, miraculin, from miracle fruit. *J. Biol. Chem.* **263**, 11536–11539 (1988).
69. Takai, A. *et al.* Secretion of miraculin through the function of a signal peptide conserved in the Kunitz-type soybean trypsin inhibitor family. *FEBS Lett.* **587**, 1767–1772 (2013).
70. Kokiladevi, E., Manickam, A. & Thayumanavan, B. Characterization of alpha-amylase inhibitor in *Vigna sublobata*. *Bot. Bull. Acad. Sin.* **46** (2005).
71. Manen, J.-F. *et al.* A nodulin specifically expressed in senescent nodules of winged bean is a protease inhibitor. *Plant Cell* **3**, 259–270 (1991).
72. Andrews, S. FastQC: A quality control tool for high throughput sequence data. Available from <http://www.bioinformatics.babraham.ac.uk>. (2010).
73. Grabherr, M. G. *et al.* Trinity: reconstructing a full-length transcriptome without a genome from RNA-Seq data. *Nat. Biotechnol.* **29**, 644 (2011).
74. Zerbino, D. R. & Birney, E. Velvet: algorithms for *de novo* short read assembly using de Bruijn graphs. *Genome Res.* **18**, 821–829 (2008).
75. Chevreux, B., Wetter, T. & Suhai, S. In *German conference on bioinformatics*. 45–56.
76. Mundry, M., Bornberg-Bauer, E., Sammeth, M. & Feulner, P. G. Evaluating characteristics of *de novo* assembly software on 454 transcriptome data: a simulation approach. *PLoS One* **7**, e31410 (2012).
77. Kumar, S. & Blaxter, M. L. Comparing *de novo* assemblers for 454 transcriptome data. *BMC Genomics* **11**, 571, doi: 10.1186/1471-2164-11-571 (2010).
78. Conesa, A. *et al.* Blast2GO: a universal tool for annotation, visualization and analysis in functional genomics research. *Bioinformatics* **21**, 3674–3676, doi: 10.1093/bioinformatics/bti610 (2005).
79. Da Maia, L. C. *et al.* SSR Locator: Tool for simple sequence repeat discovery integrated with primer design and PCR simulation. *Int. J. Plant Genomics* (2008).
80. McGinnis, S. & Madden, T. L. BLAST: at the core of a powerful and diverse set of sequence analysis tools. *Nucleic Acids Res.* **32**, W20–W25 (2004).
81. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
82. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).
83. Darrriba, D., Taboada, G. L., Doallo, R. & Posada, D. ProtTest 3: fast selection of best-fit models of protein evolution. *Bioinformatics* **27**, 1164–1165 (2011).
84. Ronquist, F. *et al.* MrBayes 3.2: efficient Bayesian phylogenetic inference and model choice across a large model space. *Syst. Biol.* **61**, 539–542 (2012).
85. Rambaut, A., Suchard, M. A., Xie, D. & Drummond, A. J. Tracer v1.5, available from <http://beast.bio.ed.ac.uk/Tracer>. (2013).

## Acknowledgements

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

## Author Contributions

All authors contributed to various aspects of this work (ordered by degree of contribution): conceived the study (A.N.E. and N.S.), aided in study design (A.N.E., N.S., M.V. and P.S.), obtained funds for the research (J.J.D. and A.N.E.), coordinated activities (A.N.E. and N.S.), obtained and grew plants from seed (A.N.E.), extracted RNA and prepared cDNA libraries for 454 sequencing (A.N.E.), conducted bioinformatic analyses (M.V., P.S., A.N.E. and R.C.), and contributed to preparation of the manuscript (M.V., A.N.E., N.S., P.S., J.J.D. and R.C.). All authors reviewed the manuscript.

## Additional Information

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

**Supplementary information** accompanies this paper at <http://www.nature.com/srep>

**Competing financial interests:** The authors declare no competing financial interests.

**How to cite this article:** Vatanparast, M. *et al.* Transcriptome sequencing and marker development in winged bean (*Psophocarpus tetragonolobus*; Leguminosae). *Sci. Rep.* **6**, 29070; doi: 10.1038/srep29070 (2016).



This work is licensed under a Creative Commons Attribution 4.0 International License. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in the credit line; if the material is not included under the Creative Commons license, users will need to obtain permission from the license holder to reproduce the material. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>