

RESEARCH ARTICLE

Open Access

De novo transcriptome assembly of *Pueraria montana* var. *lobata* and *Neustanthus phaseoloides* for the development of eSSR and SNP markers: narrowing the US origin(s) of the invasive kudzu

Matthew S. Haynsen^{1,2,3}, Mohammad Vatanparast³, Gouri Mahadwar^{3,4}, Dennis Zhu^{3,5}, Roy Z. Moger-Reischer^{3,6}, Jeff J. Doyle⁷, Keith A. Crandall^{2,8} and Ashley N. Egan^{3*}

Abstract

Background: Kudzu, *Pueraria montana* var. *lobata*, is a woody vine native to Southeast Asia that has been introduced globally for cattle forage and erosion control. The vine is highly invasive in its introduced areas, including the southeastern US. Modern molecular marker resources are limited for the species, despite its importance. Transcriptomes for *P. montana* var. *lobata* and a second phaseoloid legume taxon previously ascribed to genus *Pueraria*, *Neustanthus phaseoloides*, were generated and mined for microsatellites and single nucleotide polymorphisms.

Results: Roche 454 sequencing of *P. montana* var. *lobata* and *N. phaseoloides* transcriptomes produced read numbers ranging from ~ 280,000 to ~ 420,000. Trinity assemblies produced an average of 17,491 contigs with mean lengths ranging from 639 bp to 994 bp. Transcriptome completeness, according to BUSCO, ranged between 64 and 77%. After vetting for primer design, there were 1646 expressed simple sequence repeats (eSSRs) identified in *P. montana* var. *lobata* and 1459 in *N. phaseoloides*. From these eSSRs, 17 identical primer pairs, representing inter-generic phaseoloid eSSRs, were created. Additionally, 13 primer pairs specific to *P. montana* var. *lobata* were also created. From these 30 primer pairs, a final set of seven primer pairs were used on 68 individuals of *P. montana* var. *lobata* for characterization across the US, China, and Japan. The populations exhibited from 20 to 43 alleles across the seven loci. We also conducted pairwise tests for high-confidence SNP discovery from the kudzu transcriptomes we sequenced and two previously sequenced *P. montana* var. *lobata* transcriptomes. Pairwise comparisons between *P. montana* var. *lobata* ranged from 358 to 24,475 SNPs, while comparisons between *P. montana* var. *lobata* and *N. phaseoloides* ranged from 5185 to 30,143 SNPs.

(Continued on next page)

* Correspondence: egana@si.edu; ashegan2@gmail.com

³Department of Botany, National Museum of Natural History, Smithsonian Institution, Washington, DC, USA

Full list of author information is available at the end of the article



(Continued from previous page)

Conclusions: The discovered molecular markers for kudzu provide a starting point for comparative genetic studies within phaseoloid legumes. This study both adds to the current genetic resources and presents the first available genomic resources for the invasive kudzu vine. Additionally, this study is the first to provide molecular evidence to support the hypothesis of Japan as a source of US kudzu and begins to narrow the origin of US kudzu to the central Japanese island of Honshu.

Keywords: *Pueraria montana* var. *lobata*, Kudzu, *Neustanthus phaseoloides*, Transcriptome, Invasive, Molecular markers

Background

Pueraria montana (Lour.) Merr. var. *lobata* (Willd.) Maesen & Almeida ex Sanjappa and Pradeep (kudzu) and *Neustanthus phaseoloides* (Roxburgh) Benthham (tropical kudzu), members of the phaseoloid clade of subfamily Papilionoideae of the Fabaceae family, are twining vines native to Southeast Asia that have been introduced globally for livestock forage, nitrogen soil enrichment, and erosion control [1]. Prior to recent molecular and taxonomic revision [2], *Neustanthus* was placed within *Pueraria*, along with ~17 additional species native to southeast Asia [3]. A comprehensive molecular systematic study of *Pueraria* sensu van der Maesen [4] confirmed that its species, including several legumes of economic importance, comprise a polyphyletic assemblage of separate evolutionary lineages spread across the phaseoloid clade [5].

Both kudzu and tropical kudzu share a penchant for invasiveness in their naturalized areas, the southeastern United States (US) and the pantropics, respectively. Of the two taxa, kudzu is a far greater agricultural pest and has garnered the majority of scientific inquiry. Kudzu was introduced into the US during the Centennial Exposition of 1876 in Philadelphia, Pennsylvania [6]. The vine is currently found in 30 states and is considered an agricultural pest throughout the southeastern US [7], costing millions of dollars in eradication and management measures annually [8, 9]. A major aspect that could be influencing the invasiveness and spread of kudzu are high levels of genetic variation observed across populations in the US. This could be due to multiple introductions from its native range, either of a single genetically diverse population, or from multiple genetically distinct subpopulations, potentially from different geographic regions or from more than one of the taxonomically recognized varieties of *Pueraria montana*.

Several molecular markers have been used over the past two decades to estimate the introduced and native genetic diversities of kudzu and two other *Pueraria montana* varieties: *Pueraria montana* var. *montana* and *Pueraria montana* var. *thomsonii* (Benth.) Wiersema ex D.B. Ward [10–15]. However, despite the ecological and economic importance of kudzu, its modern molecular marker resources are

limited, lagging particularly in the characterization and development of microsatellites (SSRs) and single nucleotide polymorphisms (SNPs). Transcriptome sequencing is currently one of the most popular applications of next-generation sequencing due to its versatility, cost efficiency, and suitability for use on non-model organisms [16]. Transcriptomes are often mined for expressed simple sequence repeats (eSSRs) for marker development and genetic diversity studies. eSSRs have been shown to have greater transferability across taxa than traditional ‘anonymous’ SSRs [17, 18]. This increased transferability can be utilized in multiple ways. First, if a transcriptome is not available for the species of interest, a closely related species whose transcriptome is available can be used as a surrogate reference for microsatellite development. Second, if a researcher is studying two closely related taxa and transcriptomes are available for both, a single set of markers can be developed that work on both species to reduce costs. To this end, we have compared the transcriptomes of kudzu and tropical kudzu to identify shared eSSRs between the species in order to develop primers that can be used equally well for population genetic studies of either species, and shed light on the introduction history of the notorious invasive kudzu in the United States.

In the present study, three transcriptomes, two *P. montana* var. *lobata* and one *N. phaseoloides*, were de novo assembled and characterized. Intra- and inter-specific comparisons were made between transcriptomes and two sets of population genetic markers were identified, eSSRs and SNPs. The eSSRs were validated across Asian and North American populations of *P. montana*. Var. *lobata* and used to explore population diversity and structure across native and introduced ranges. The resulting data provide genetic resources for future studies of kudzu and related genera through development of high-resolution marker sets for genetic diversity assessment and population studies.

Results

Transcriptome sequencing and quality control

Transcriptome sequencing produced between 279,109 and 423,426 reads per transcriptome (Table 1), with *Neustanthus phaseoloides* (hereafter CPP02) having the

t1.1 **Table 1** Statistics following ConDeTri cleaning and Trinity assembly

t1.2	Accessions	CPP27	Pmnk6	CPP02
t1.3	Number of raw reads	279,109	396,022	423,426
t1.4	Number of raw bases (bp)	112,337,841	247,596,818	158,214,933
t1.5	Number of clean reads	257,015	381,166	348,529
t1.6	Cleaned reads / Raw reads (%)	92.1%	71.0%	82.3%
t1.7	Number of clean bases (bp)	75,672,645	124,810,371	87,666,889
t1.8	Mean clean read length (bp)	294	444	252
t1.9	Number of aligned reads	99,248	116,524	119,452
t1.10	Aligned read / Cleaned reads (%)	38.6%	41.4%	34.3%
t1.11	Number of contigs	18,325	15,736	18,412
t1.12	Number of bases in contigs (bp)	11,703,977	15,640,762	11,892,992
t1.13	Mean contig length (bp)	639	994	646
t1.14	N50 (bp)	755	1256	759
t1.15	Longest contig (bp)	4335	4815	6221
t1.16	Number of singletons	60,869	45,306	73,994
t1.17	Singletons / Cleaned reads (%)	23.7%	16.1%	21.2%
t1.18	Number of bases in singletons (bp)	17,591,281	20,431,176	18,048,611
t1.19	Mean singleton length (bp)	289	451	244
t1.20	Number of transcripts (contigs + singletons)	79,194	61,042	92,406
t1.21	<i>bp</i> base pairs			

128 most reads produced. CPP02 and the greenhouse-raised
 129 kudzu (hereafter CPP27) were sequenced on the same
 130 run and were multiplexed with two other transcriptomes
 131 not reported here. While sequencing of CPP02 produced
 132 the most reads, the mean read length before cleaning
 133 was shorter than that of CPP27 (373 bp vs. 402 bp, re-
 134 spectively), as was the mean read length after cleaning
 135 (252 vs. 294, respectively). The tendency for shorter
 136 DNA fragments to be incorporated at the library con-
 137 struction phase and sequencing stage may provide an
 138 explanation for the difference in the number of raw
 139 reads produced between CPP27 and CPP02. However,
 140 following cleaning, the number of clean bases was com-
 141 parable between CPP02 and CPP27, as were all other
 142 downstream metrics (Table 1). While 454 pyrosequenc-
 143 ing was used for all three transcriptomes, the chemis-
 144 tries between the two CPP transcriptomes and the
 145 wild-collected kudzu (hereafter Pmnk6) transcriptome
 146 differed, with the Pmnk6 transcriptome benefiting from
 147 an improved chemistry, as seen in the increased number
 148 of raw bases, the average read length before cleaning
 149 (625 bp) and the mean clean read length (444; Table 1).
 150 These sequencing improvements translated into im-
 151 proved assembly statistics, such as increased mean con-
 152 tig length (~1.5× that of the CPP transcriptomes),
 153 higher N50 (1.65× CPP) and fewer singletons (Table 1).
 154 However, the improved chemistry did not lead to differ-
 155 ences in the number of aligned reads in the assembled
 156 transcriptomes (Additional file 1).

De novo assembly

Trinity used an average of 38.1% of the ConDeTri cleaned
 reads in its assemblies and produced an average of 17,491
 contigs. The mean contig lengths ranged from 639 bp to
 994 bp (Table 1) and each of the accessions had contigs
 exceeding 3000 bp (Fig. 1). Additionally, Bowtie2 mapped
 ~68% of each accession's contigs back to their raw reads
 (Additional file 1). Overall transcriptome contamination
 was low, with fungal contamination ranging between 2.64
 and 3.53%, while prokaryote and viral contamination
 ranged from 0.5 to 1.32% (Additional file 2). Transcrip-
 tome completeness varied greatly, with a range of
 complete units from 164 to 361 and duplicate units simi-
 larly showing a >2× difference between transcriptomes
 (Fig. 2). Specifically, transcriptome completeness was ap-
 proximately 64, 77, and 70%, for CPP27, Pmnk6, and
 CPP02, respectively. The reciprocal best BLAST hits
 (RBH) of the transcriptomes showed that 1525 transcripts
 were shared among all three (Fig. 3).

Functional annotation of transcriptomes

In total, we have obtained 13,230, 18,446 and 24,447 asso-
 ciated GO IDs for CPP02, CPP27 and Pmnk6 transcrip-
 tomes, respectively (Table 2) corresponding to the 33, 43
 and 51% of original contigs in each transcriptome, while
 only 9.6, 17 and 36% of the singletons had associated func-
 tional protein information (GO IDs). Therefore, more
 than 90, 82 and 63% of singletons were discarded during
 the multiple searches, which is unfortunate because over

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

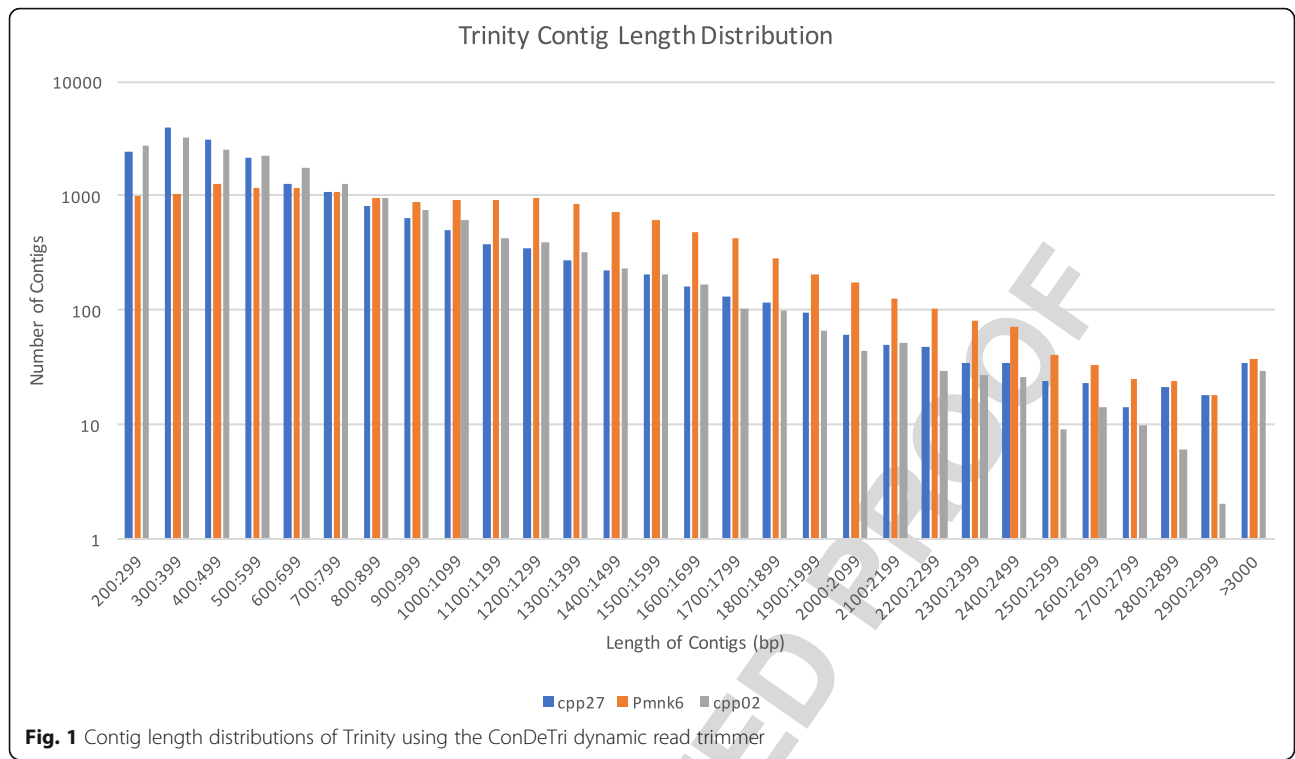
181

182

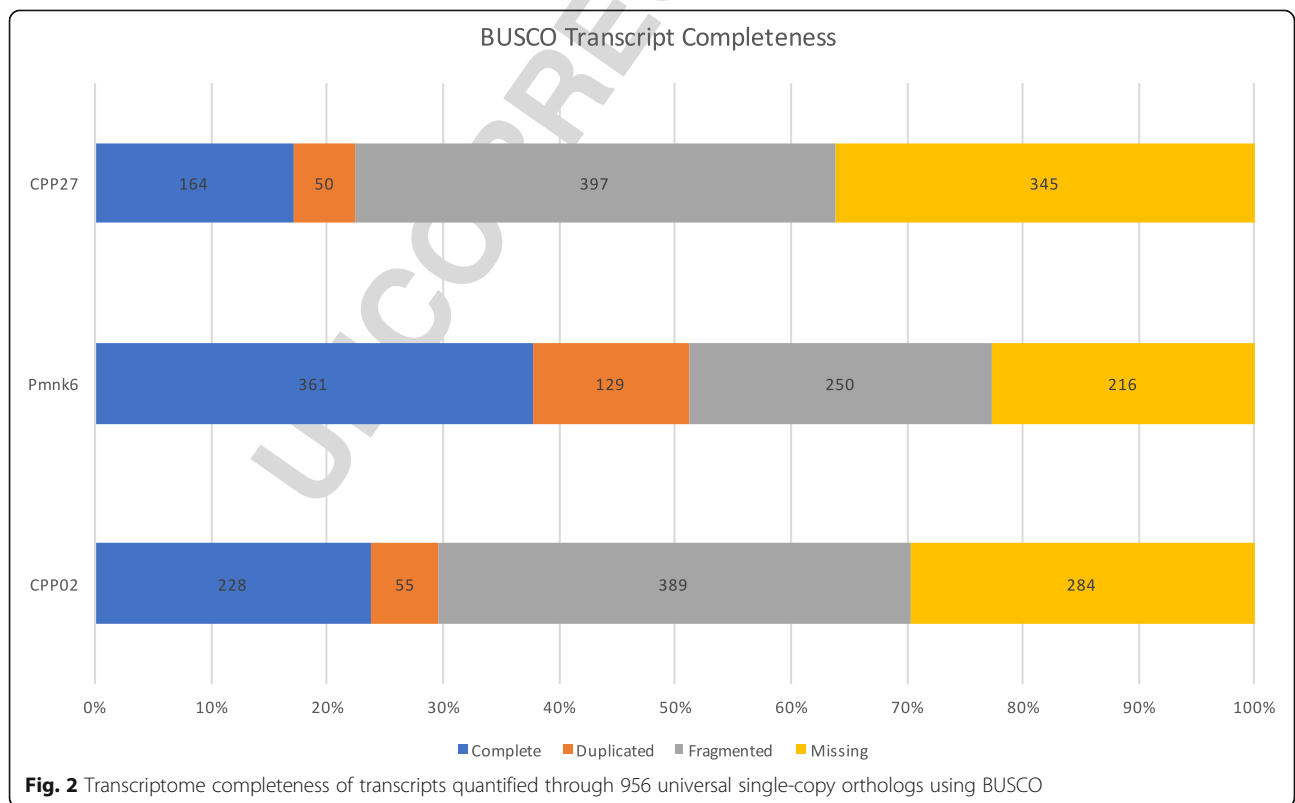
183

184

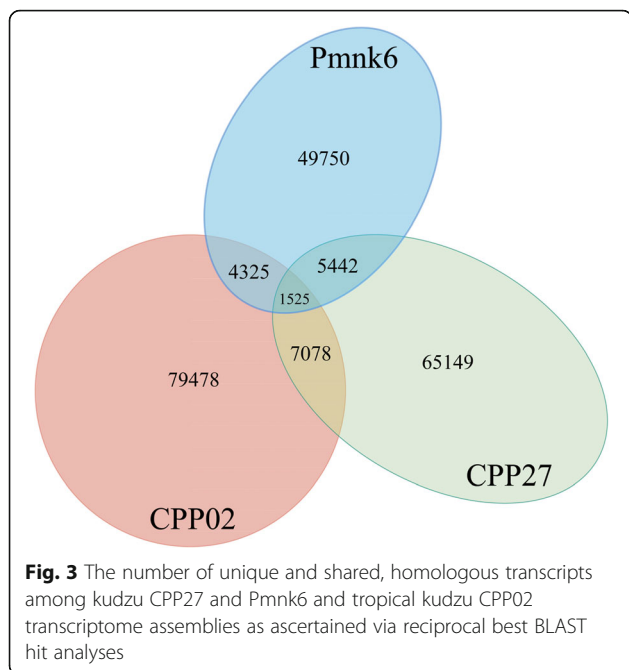
184



f1.1
f1.2



f2.1
f2.2



f3.1
f3.2
f3.3
f3.4
f3.5

tropical kudzu (CPP27 vs. CPP02, Pmnk6 vs. CPP02, and CPP27/Pmnk6 vs. CPP02, respectively). The over 30,000 SNPs identified between CPP27/Pmnk6 vs. CPP02 is greater than the sum of SNPs from the individual comparisons of *P. montana* var. *lobata* to *N. phaseoloides* because the merged transcripts offer a more complete snapshot of a US kudzu transcriptome which was used as the reference for SNP detection. Lastly, we found 24,475 SNPs within kudzu from among three countries (Japan vs. Pmnk6(US)/CPP27(US)/ China). The majority of high-confidence SNPs were found within contigs rather than singletons (Table 3), which is expected given the fact that more highly expressed genes will be more likely to be represented by >20x coverage (one of our criteria for high confidence) and are most likely to assemble into contigs. Also of note, the transition/transversion ratio varied from 1.41 to 1.73 (Table 3) with the higher ratios found between the intergeneric comparisons than the intraspecific comparisons.

eSSR discovery and characterization

The eSSR analysis of the transcripts detected 5255 and 4586 perfect eSSRs for CPP27 and CPP02, respectively. The majority (76.7 and 76.8%) of eSSRs were tri-nucleotide repeats (TNRs; Table 4). After vetting for primer design, there were 1646 potential eSSRs identified in *P. montana* var. *lobata* and 1459 in *N. phaseoloides*. Looking only at TNRs (1458 for CPP27 and 1273 for CPP02), 25 matches were found between *P. montana* var. *lobata* and *N. phaseoloides* in which either the forward or reverse primers were identical, suggesting homology. However, no sets of primer pairs (forward and reverse primers together) were found duplicated between transcriptomes. Alterations to the non-identical primer pair within the 25 matches allowed for the creation of 17 identical primer pairs between CPP27 and CPP02. These 17 shared primer pairs represent inter-generic phaseoloid eSSRs. Additionally, 13 TNR primer pairs specific to *P. montana* var. *lobata* were also selected for screening. Of the 30 total eSSR primer pairs, 21 pairs were advanced to the Culley et al. [19] protocol; of the nine primer pairs that were eliminated, four did not amplify a product, four amplified in an unexpected

T4

185 54, 56 and 66% of final annotated transcripts belong to
186 the singletons in CPP02, CPP27, and Pmnk6, respectively
187 (Table 2). In all three transcriptomes, the highest top hit
188 species for the annotated proteins were *Glycine max* (L.)
189 Merr., *G. soja* Siebold & Zucc. and *Cajanus cajan* (L.)
190 Millsp., respectively (Additional files 3, 4, and 5). Summar-
191 ies of the biological process, cellular components and mo-
192 lecular function categories for each transcriptome are
F4 193 shown in Fig. 4.

194 **SNP discovery**
195 We conducted pairwise tests for high-confidence SNP
T3 196 discovery of the kudzu transcriptomes (Table 3, Add-
197 itional files 6, 7, 8, 9 and 10). Our conservative assess-
198 ment of SNPs reduced thousands of high-confidence
199 SNPs to a lower number (Table 3) that are 1) one-to-one
200 point mutations without length variants, 2) have vari-
201 ation frequency over 95%, and 3) have a repeat depth of
202 20 or more. As such, we identified 358 SNPs between
203 the two US kudzu transcriptomes (CPP27 vs. Pmnk6),
204 and 5185, 19,028, and 30,143 SNPs between kudzu and

t2.1 **Table 2** Summary of gene ontology analysis

t2.2	Accessions	Transcripts	Orfs	Predictions	BLAST Hits	Annotated GO IDs	ECs
t2.3	CPP27	79,194	37,741	30,716	28,795	18,446	8039
t2.4		(18,325/60869)		(13,534/17182)	(12,583/16212)	(7958/10488)	
t2.5	Pmnk6	61,042	50,320	42,386	39,366	24,447	6337
t2.6		(15,736/45306)		(14,821/27565)	(12,705/26661)	(8079/16368)	
t2.7	CPP02	92,406	34,223	27,661	22,472	13,230	4064
t2.8		(18,412/73994)		(14,677/12984)	(10,407/12065)	(6085/7145)	

t2.9 *Orfs* open reading frames, *GO* gene ontology, *ECs* enzyme codes. Parentheses: (contigs/singletons)

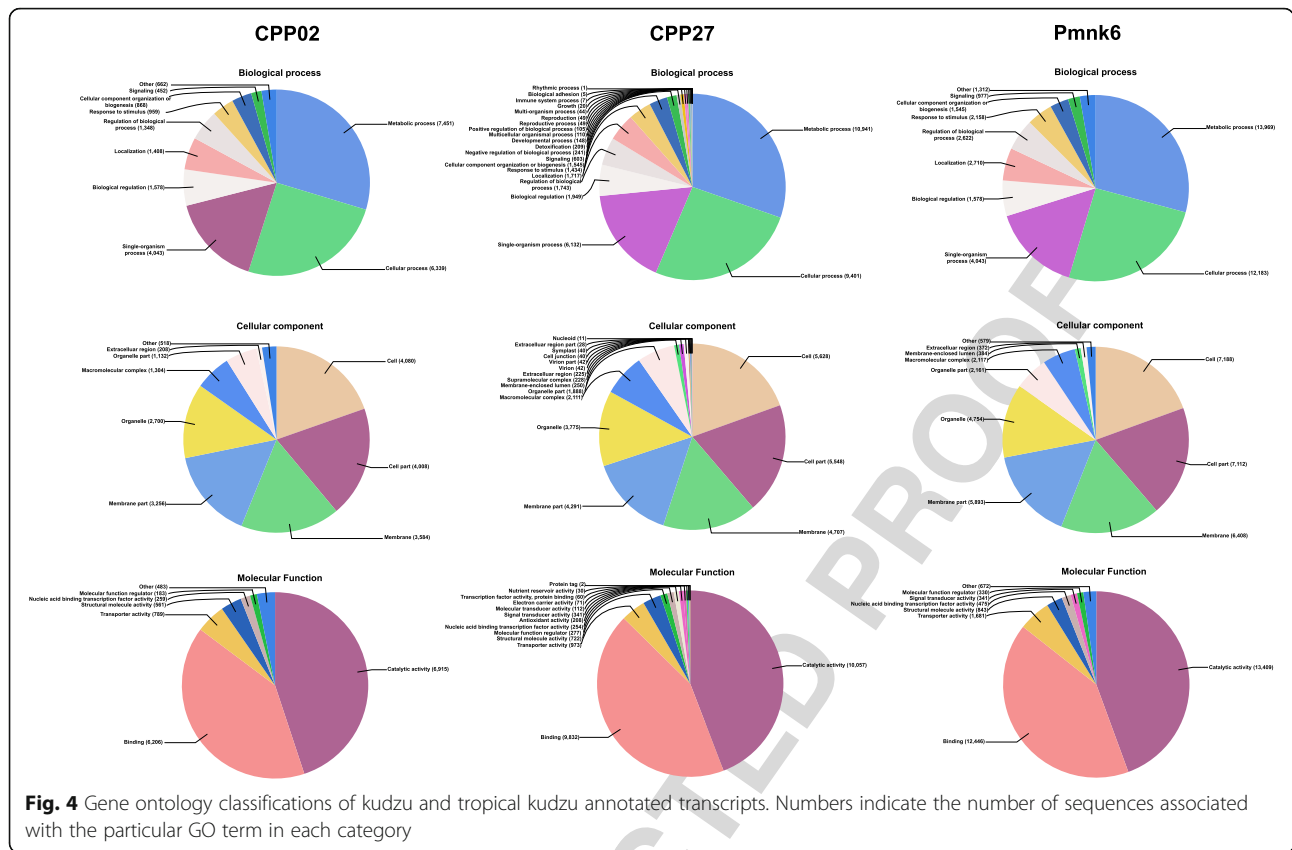


Fig. 4 Gene ontology classifications of kudzu and tropical kudzu annotated transcripts. Numbers indicate the number of sequences associated with the particular GO term in each category

f4.1
f4.2
f4.3

247 size range, and one displayed double banding (Addi-
248 tional file 11). Of the 21 primer pairs that were assessed
249 with the Culley et al. [19] protocol, seven were discarded
250 due to multiple banding and four for lack of amplifica-
251 tion, whereas a further three were removed due to the
252 presence of monomorphic alleles (Additional file 11).
253 The final set of eSSR primer pairs identified seven poly-
254 morphic loci displaying single bands of expected sizes
T5 255 (Table 5).

Population structure and genetic diversity of kudzu

256
257 Three genetic units were determined to be the optimal
258 value of K in STRUCTURE across the 75 accessions (K
259 = 3, Fig. 5, Additional file 12). The US is primarily com-
260 posed of a single genetic unit, with a couple individuals
261 assigned to a second unit; whereas, China and Japan are
262 more heterogeneous in their composition, yet they are
263 still composed of the same 2 units found in the US.
264 Thailand, however, is composed of a single genetic unit

F5

Table 3 Single nucleotide polymorphism detection among kudzu and tropical kudzu genotypes

Comparison	HC SNPs	SNPs > 95% ^a	SNPs > 20x ^b	Total SNPs ^c	Ts/Tv
Pmnk6 vs CPP27	10,417 (7494/2923)	6016 (4125/1891)	426 (252/174)	358	1.41
CPP02 vs CPP27	99,584 (81,276/18308)	86,626 (70,638/15988)	5831 (5091/740)	5185	1.60
CPP02 vs Pmnk6	220,739 (164,118/56621)	174,884 (127,311/47573)	21,258 (19,255/2003)	19,028	1.73
CPP02 vs Pmnk6, CPP27	314,416 (229,163/85251)	248,719 (178,102/70617)	33,603 (29,812/3791)	30,143	1.71
Japan vs Pmnk6, CPP27, China	494,234 (494,234/0)	79,088 (79,088/0)	27,108 (27,108/0)	24,475	1.47

t3.13 ^aSNPs with the > 95% frequency
t3.14 ^bSNPs with > 95% frequency and > 20x coverage
t3.15 ^cOne-to-one point mutations after exclusion of indels and length variants; HC: high confidence; parentheses: (contigs/singletons)

Table 4 Transcriptome eSSRs

	CPP27	CPP02
t4.2		
t4.3	79,194	92,406
t4.4	5255	4586
t4.5	770	670
t4.6	4032	3524
t4.7	180	137
t4.8	106	79
t4.9	167	176
t4.10	1646	1459
t4.11	14	28
t4.12	1458	1273
t4.13	62	54
t4.14	41	25
t4.15	71	79

sampled. Genetic structuring as assessed by pairwise F_{st} showed differences among groups, particularly in Thailand and southern China (China 3; Table 8), corroborated by the structuring of genetic units shown in Fig. 5. As defined by Wright [20], Thailand showed very great genetic variation ($F_{st} > 0.25$) with respect to all other subpopulations, except China 3, with which it showed great variation ($0.15 < F_{st} < 0.25$). The rest of the comparisons resulted in little to moderate genetic variation ($0 < F_{st} < 0.05$ and $0.05 < F_{st} < 0.15$, respectively). The neighbor-joining distance tree supports the pairwise F_{st} results (Fig. 6): 1) Thailand is a distantly related lineage to the nine other subpopulations representing *P. montana* var. *montana* and var. *thomsonii*; 2) the Chinese subpopulations are divided into three lineages; and 3) the US subpopulations are more genetically similar to Japan 2.

Discussion

Invasive species are increasingly widening their scope across the globe, yet the genetic mechanisms underlying invasiveness or weediness remain a mystery. In the genomics era, scientists have raised a clarion call to arms to build genomic resources to study invasive species [21]. Understanding the introduction history and relative genetic diversity of invasive species is an important step to gaining a foothold on management and control, a goal requiring the development of variable molecular markers such as microsatellites or SNPs to assess genetic diversity and population structure. In this study, we have assembled and characterized multiple transcriptomes of the invasive Kudzu vine, *Pueraria montana* var. *lobata*, and for tropical kudzu, *Neustanthus phaseoloides*, a species until recently thought congeneric with kudzu [2, 5].

that is unique to that nation, which supports our classification of its accessions as being different varieties of *P. montana*, specifically var. *thomsonii* and var. *montana*.

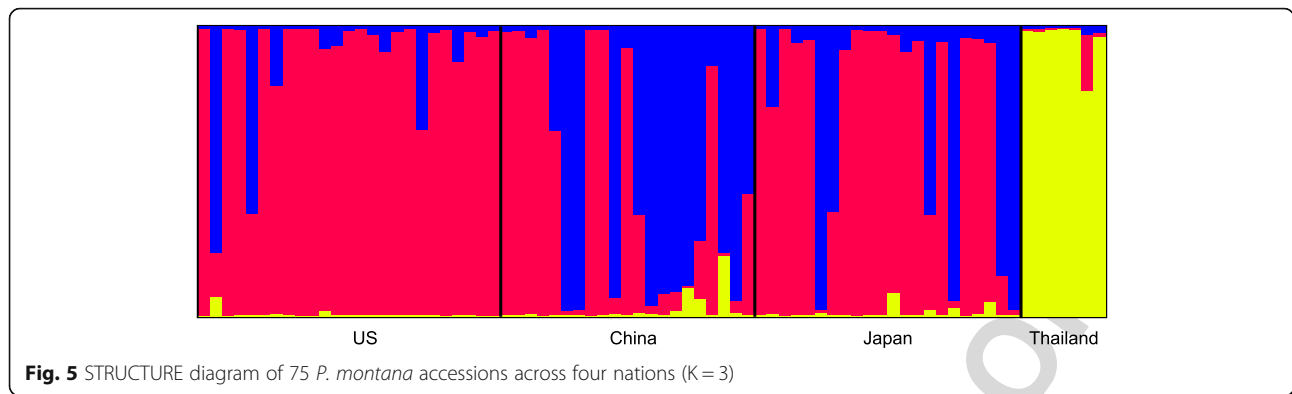
The national populations exhibited from 20 to 43 alleles across a total of seven loci (Table 6), while the subpopulations exhibited from 20 to 36 total alleles (Additional file 13). China was composed of the greatest number of alleles, in particular, China 3 (southern), while Thailand was composed of the fewest number of alleles.

After Bonferroni correction, none of the subpopulations' observed and expected heterozygosities significantly differed (Table 7), supporting the hypothesis that all the subpopulations were in Hardy-Weinberg equilibrium when

Table 5 Seven eSSR primers optimized and used to assess population genetics in kudzu accessions

Locus	Sequence	Dye/Tail	SSR	Length (bp)
PP2	F: 5'-TAG GAG TGC AGC AAG CAT ATG CCG CGG ATC TTT GAA AG-3'	VIC /M13A	AAC	100–130
	R: 5'-CAA ATT GGC CCT GTC CCA AT-3'	n/a		
PP4	F: 5'-TGT AAA ACG ACG GCC AGT CAT GCC CAC GTG CTT CAT AG-3'	6FAM/M13	GCT	100–140
	R: 5'-CTC TCA GAT CCA GGC CCA AA-3'	n/a		
PP10	F: 5'-TAG GAG TGC AGC AAG CAT GGC ATG TAG ATC CAG CTA AA-3'	VIC/M13A	GGT	310–330
	R: 5'-TTG ACA GAT TTC TGA TTC TTG G-3'	n/a		
PP13	F: 5'-TAG GAG TGC AGC AAG CAT GAT TGA GCA GGC ACG AGA AC-3'	VIC/M13A	GCT	270–300
	R: 5'-CAG TAG CAG GCA TGT GTT GG-3'	n/a		
PL1	F: 5'-CAC TGC TTA GAG CGA TGC TGT AAG CGT TCG TTC GTT GG-3'	PET/M13B	CTT	400–440
	R: 5'-TCA ACC TGG TGC TCT CTG AC-3'	n/a		
PL7	F: 5'-TGT AAA ACG ACG GCC AGT AGT GGC CTT GCT CTT CTT CC-3'	6FAM/M13	CTT	80–140
	R: 5'-GTG TCA TCT CAG CAC GTT GG-3'	n/a		
PL11	F: 5'-TGT AAA ACG ACG GCC AGT TGG CAT CAT CCT TCA ACC AC-3'	6FAM/M13	ACC	300–330
	R: 5'-ATT CCG GAA TAG TGG GTG GG-3'	n/a		

F forward primer, R reverse primer. Dyes VIC: 2'-chloro-7'-phenyl-1,4-dichloro-6-carboxy-fluorescein; 6FAM: 6-carboxyfluorescein; PET: chemical structure currently unpublished as proprietary to Lifetech. Tail: see Culley et al. [19] for information about M13, M13A, and M13B



f5.1
f5.2

311 Kudzu is well known as an invasive species in both agri-
 312 cultural and natural areas due to its fast growth, clonal
 313 habit, and extensive introductions outside its native range.
 314 Tropical kudzu is also known to be invasive in its intro-
 315 duced ranges, but to a lesser extent. We mined our tran-
 316 scriptomes of these two species for molecular markers
 317 (eSSRs and SNPs), screened and validated eSSRs, and per-
 318 formed functional annotations of the transcriptomes, im-
 319 proving the genetic resources available for kudzu and
 320 tropical kudzu.

321 **Transcriptome characterization**

322 Whether researching model or non-model organisms, se-
 323 quencing the transcriptome of a species is a natural begin-
 324 ning for genome-wide resource development and study
 325 [22, 23], enabling the characterization of gene expression
 326 profiles, genetic marker discovery, and phylogenetic infer-
 327 ence [24]. Here, we characterize the transcriptomes of two
 328 accessions of kudzu, one wild-collected (Pmnk6) and one
 329 partially inbred line propagated by the USDA agriculture
 330 research service (CPP27), as well as one of tropical kudzu
 331 (CPP02). We chose to use 454 pyrosequencing technology
 332 over Illumina due to the longer read lengths, an important

333 consideration when dealing with potentially polyploid 333
 334 plants [23, 25, 26]. *Pueraria* is descendent from an ancient 334
 335 polyploidy event that transpired 50–60 mya near the origi- 335
 336 n of the papilionoid subfamily [27, 28], creating a dupli- 336
 337 cated genomic complement that has fractionated over 337
 338 time but whose signature still remains within descendent 338
 339 genomes. Longer reads are more likely to unambiguously 339
 340 assemble or align across homoeologues, duplicated genes 340
 341 produced via allopolyploidy [29]. Furthermore, the longer 341
 342 reads result in the sequencing of more full-length mRNA 342
 343 transcripts, an outcome that argues for including single- 343
 344 tons (those reads that do not assemble into contigs) in the 344
 345 overall transcript complement. Although pyrosequencing 345
 346 produces fewer overall reads as compared to Illumina, its 346
 347 ability to produce longer transcripts is advantageous, particu- 347
 348 larly for allopolyploid species and other hybrids where 348
 349 avoiding the assembly of chimeric sequences is important. 349

350 The comparative results across our transcriptomes in 350
 351 terms of the number of transcripts discovered and the 351
 352 relative overlap among pairwise comparisons provides 352
 353 some insights into the relative impact of environment vs. 353
 354 shared ancestry. CPP02 had the highest number of tran- 354
 355 scripts and the highest number of unique transcripts, 355
 356 with Pmnk6 having the least number of transcripts, even 356
 357 though it presents the best transcriptome in terms of 357
 358 mean contig length, N50, and BUSCO results. One ex- 358
 359 planation involves the number of tissues used for se- 359
 360 quencing. CPP02 utilized three tissues (young leaves, 360
 361 young shoot tips, and buds) while CPP27 used two tis- 361
 362 sues (young leaves and young shoot tips), and Pmnk6 362
 363 used a single tissue (young leaves). Given this informa- 363
 364 tion, it makes sense that the transcriptome that was 364
 365 composed of the greatest number of tissues resulted in 365
 366 the highest number of unique transcripts due to expres- 366
 367 sional differences across tissue types. CPP02 and CPP27 367
 368 shared the highest number of reciprocal best BLAST hits 368
 369 (RBH). However, one would expect the two kudzu acces- 369
 370 sions (CPP27 and Pmnk6) to share the greatest number of 370
 371 overlapping transcripts due to shared ancestry. This could 371
 372 also be explained by the fact that the two transcriptomes 372

t6.1 **Table 6** Allelic frequency for *Pueraria* national populations

t6.2 Locus	USA N = 25	China N = 21	Japan N = 22	Thailand N = 7	Mean	SD	Total
t6.3 PP2	8	7	6	4	6.25	1.71	9
t6.4 PP4	4	5	7	3	4.75	1.71	9
t6.5 PP10	5	5	6	3	4.75	1.26	8
t6.6 PP13	3	7	4	2	4.00	2.16	7
t6.7 PL1	4	4	2	4	3.50	1.00	9
t6.8 PL7	8	8	11	3	7.50	3.32	15
t6.9 PL11	5	7	3	1	4.00	2.58	7
t6.10 Mean	5.29	6.14	5.57	2.86	4.96	1.96	9.14
t6.11 SD	1.98	1.46	2.99	1.07	1.42	0.80	2.73
t6.12 Total	37	43	39	20	34.75	13.73	64

t6.13 N number of accessions, SD standard deviation

Table 7 Observed and expected heterozygosities for *Pueraria* subpopulations

	US 1	US 2	US 3	CN 1	CN 2	CN 3	JP 1	JP 2	JP 3	TH
# Individuals	8	10	7	5	8	8	7	8	7	7
Obs. Het.	0.717	0.552	0.472	0.611	0.378	0.632	0.396	0.506	0.656	0.594
Exp. Het.	0.643	0.503	0.547	0.648	0.589	0.763	0.579	0.572	0.661	0.583
HWE p-value	0.766	0.251	0.611	0.765	0.079	0.392	0.013	0.429	0.869	0.442

US United States, CN China, JP Japan, TH Thailand, Obs: Observed, Exp Expected, Het Heterozygosity, HWE Hardy-Weinberg Equilibrium

that shared the most homologous tissues resulted in the highest number of shared transcripts. Alternatively, the seeming disparity in shared best BLAST hits could be explained by the relative impacts of a shared environment, which often affects gene expression. Our two CPP transcriptomes were both grown in the same greenhouse environment at the same time and so their gene expression profiles may be expected to be more similar than those of the two *P. montana* var. *lobata* accessions, one of which was grown in the greenhouse (CPP27) and one in the wild (Pmnk6). A similar finding was discovered across transcriptomes of *Eutrema salsugineum* (Pall.) Al-Shehbaz & Warwick plants that were grown in field (uncontrolled environment) vs. cabinet (controlled environment) conditions, with the plants grown in the controlled environment sharing a higher number of expressed genes as compared to the more geographically proximate plants grown in differing environments [30].

In this study, we were able to annotate over 13,000 transcripts from kudzu and tropical kudzu (Table 1). Our transcriptomes do not provide a full gene complement due to low sequencing depth as evidence by our BUSCO results (Fig. 2). However, the level of unannotated transcripts in this study is similar to results reported from other non-model legumes, like winged bean [31], chickpea [32], and field pea [33]. The unidentified transcripts are likely due to 1) correspondence to non-coding regions or pseudogenes, 2) short length of transcripts, or 3) coding genes that

have yet to be described, perhaps including species-specific “orphan” genes [34]. Catalytic activity, binding, metabolic and cellular processes were among the most highly represented groups regarding GO analysis (Fig. 4) across all three transcriptomes, as expected given that we used young tissues that are undergoing extensive metabolic activities.

Single nucleotide polymorphism discovery

SNPs are fast becoming the marker of choice due to their ease of discovery via next generation sequencing technologies [35]. Additionally, the ease of mining SNPs from previously produced transcriptomes can provide a new use for previously published data sets that may be sitting idle in online repositories. SNPs, though less polymorphic than SSRs, may provide higher resolution assessment of genetic variation and identification of population structure [36]. We detected a near 100-fold increase in the number of SNPs detected between kudzu and tropical kudzu as compared to that detected within kudzu. SNPs discovered between kudzu and tropical kudzu may represent species level, fixed differences between these genera. Validation of these SNPs is beyond the scope of this paper; nevertheless, this list presents a significant resource for future work in genetic diversity assessment, genetic mapping, genome-wide association mapping, or evolution-based studies of invasiveness, and marks the first SNP markers discovered to date in *Pueraria* and *Neustanthus*. Use of these SNP markers across

Table 8 Subpopulation pairwise F_{st}

	US 1	US 2	US 3	CN 1	CN 2	CN 3	JP 1	JP 2	JP 3	TH
US 1	–	0.811	0.541	0.441	0.009	0.000*	0.297	0.126	0.099	0.000*
US 2	–0.023	–	0.378	0.730	0.009	0.000*	0.432	0.108	0.072	0.000*
US 3	–0.011	–0.008	–	0.306	0.009	0.000*	0.360	0.946	0.153	0.000*
CN 1	–0.009	–0.022	0.024	–	0.108	0.009	0.901	0.162	0.108	0.000*
CN 2	0.075	0.098	0.099	0.075	–	0.297	0.207	0.081	0.739	0.000*
CN 3	0.077	0.107	0.120	0.073	0.022	–	0.063	0.000*	0.324	0.000*
JP 1	0.015	–0.002	0.022	–0.035	0.051	0.064	–	0.207	0.486	0.000*
JP 2	0.016	0.025	–0.030	0.049	0.078	0.085	0.042	–	0.135	0.000*
JP 3	0.029	0.028	0.042	0.037	–0.014	0.006	0.002	0.036	–	0.000*
TH	0.315	0.370	0.377	0.330	0.310	0.244	0.347	0.352	0.322	–

Below diagonal pairwise F_{st} values, above diagonal p -values

US United States, CN China, JP Japan, TH Thailand

* = significant under Bonferroni correction ($p < 0.001$)

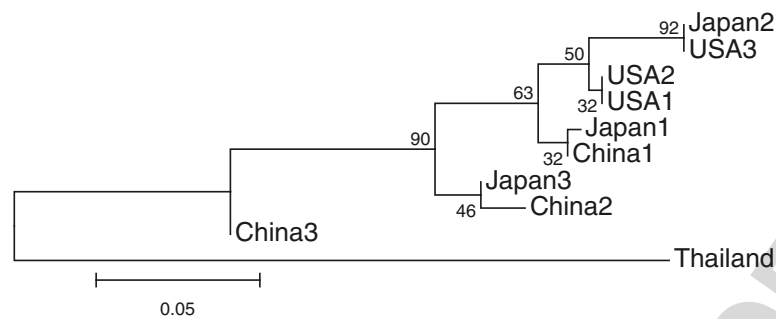


Fig. 6 Neighbor joining distance tree based on F_{st} values and 10,000 bootstraps. US = United States; CN = China; JP = Japan; and TH = Thailand

f6.1
f6.2

428 a wide population-level sampling throughout Asia would
429 enable a robust investigation into the introduction his-
430 tory of kudzu within the US.

431 eSSR marker discovery and validation

432 eSSRs are routinely developed from transcriptomic data,
433 providing a ready source for genetic diversity assessment
434 through cost-effective means [37]. In spite of being derived
435 from coding DNA, which is evolutionarily conserved,
436 eSSRs have proven a variable and valuable resource for gen-
437 etic studies [18]. In our study, we detected ~5000 eSSRs
438 each within kudzu and tropical kudzu. Overall, trinucleo-
439 tide SSR motifs (TNRs) were the most abundant, as found
440 consistently in other plant studies [17, 38–42]. Presumably
441 this is because TNRs will not affect the open reading frames
442 of coding regions [38]. We investigated the utility of 30
443 eSSR markers discovered in our data and optimized seven
444 for use across kudzu. When compared to the
445 kudzu-derived SSR markers of Hoffberg et al. [14], similar-
446 ities and benefits are found. For instance, Hoffberg et al.
447 [14] assessed their 15 genomic SSRs against 102 geograph-
448 ically dispersed individuals, finding that their alleles per
449 locus ranged from 2 to 8, whereas our alleles per locus
450 ranged from 7 to 15 (Table 6). This comparison shows
451 twice as many alleles within a smaller sample size, approxi-
452 mately two-thirds the size of Hoffberg et al. [14]. One ex-
453 planation for the difference in allele numbers could be
454 attributed to the differing sampling ranges, with our indi-
455 viduals being collected from a greater global area. However,
456 when Bentley and Mauricio [15] used the Hoffberg et al.
457 [14] primers on 1747 accessions of kudzu from solely the
458 US they identified 2–17 alleles per locus, which also repre-
459 sents a doubling of alleles but in a smaller sampling area.
460 Additionally, when our observed heterozygosities are com-
461 pared to the primers of Hoffberg et al. [14], they ranged
462 from 0.372–0.726 (Table 7), while Hoffberg et al. [14]
463 ranged from 0.0–0.9 and Bentley and Mauricio ranged from
464 0.004–0.741. The large difference in the heterozygosity
465 comparisons, particularly when focusing on the low end,
466 may be attributed to differences in sampling strategies.
467 Bentley and Mauricio [15] report sampling kudzu within a

468 population every few meters, suggesting that they treated a
469 patch of kudzu as a population, whereas we sampled indi-
470 viduals no closer than ~1 km apart, and viewed a popula-
471 tion as a regional area comprised of numerous,
472 non-connected patches. With the abilities to grow over
473 12 in. per day and root at the nodes, a kudzu patch may
474 likely represent only one or a few genets [43]. Therefore,
475 the reported clonal sampling of Bentley and Mauricio [15]
476 may be the cause of the near 0.0 observed heterozygosities
477 and may not be indicative of the primers themselves.

478 Genetic diversity of kudzu

479 For the past two decades, the genetic diversity of kudzu
480 has been assessed with the various molecular markers of
481 the corresponding era. For instance, Pappert et al. [10]
482 used 13 allozymes across 1000 US accessions to con-
483 clude that introduced kudzu possessed considerable gen-
484 etic variation with a lack of geographic structuring.
485 Similar conclusions were subsequently reached by Jewett
486 et al. [11] using 18 random amplified polymorphic DNA
487 (RAPD) markers across 50 accessions from the US and
488 China, and by Sun et al. [12] using 11 inter-simple se-
489 quence repeat (ISSR) markers across 108 accessions
490 from the US and China. A decade later, Bentley and
491 Mauricio [15], using 15 SSRs and one chloroplast
492 marker across 1747 US accessions, reported that the
493 high levels of genetic diversity result from high clonal
494 reproduction in kudzu, as described by Ellstrand and
495 Roose [44], Balloux et al. [45], and Halkett et al. [46].
496 Specifically, high levels of genetic variation are expected
497 in clonal populations when the populations were
498 founded by sexual propagules [44], which can be the
499 case even if recruitment of sexual offspring into estab-
500 lished populations is rare. This may be the case for
501 kudzu due to its deliberate introduction by landowners
502 into novel habitats from seed stock. Additionally, clonal
503 populations are capable of maintaining higher genetic
504 diversity at each locus even though they support a lower
505 number of different genotypes [45, 46]. Our results cor-
506 roborate the findings that introduced kudzu displays
507 high levels of genetic variation throughout the US (Table

508 6, Additional file 13); however, we still maintain that the
509 high genetic variation is possibly indicative of multiple
510 introductions from across its native range.

511 Population structure and introduction history of kudzu

512 Kudzu is said to have first been brought to the US by the
513 Japanese who planted it as an ornamental vine outside
514 their pavilion at the 1876 World's Fair Centennial Exhib-
515 ition in Philadelphia [47]. Later, David Fairchild, a plant
516 explorer for the United States Department of Agriculture,
517 noted its uses, including as forage, in Japan and brought
518 back some seeds to plant near his home in Washington,
519 D.C., as a trial. In the 1930's, the US government began
520 planting millions of seedlings across the southeastern
521 states as a means of erosion control. Whether the US gov-
522 ernment sourced these kudzu seedlings from one or mul-
523 tiple native populations from Japan or elsewhere is not
524 known.

525 Although there is consensus across most studies show-
526 ing robust findings of high levels of genetic variation of
527 kudzu in the US, most of the studies reported a lack of
528 geographic patterning of genotypes, and none included
529 wide sampling across Asia so as to enable an investiga-
530 tion into source populations of US introduction(s). Our
531 results include new clues in identifying the native origins
532 of US kudzu. The Thailand subpopulation is composed
533 of non-*P. montana* var. *lobata* individuals. With evi-
534 dence for strong genetic differentiation and zero popula-
535 tion admixture between Thailand and other
536 subpopulations, we can definitively rule Thailand out as
537 a source of US kudzu introductions. It may also be pos-
538 sible to rule southern China out as an origin of US
539 kudzu introductions due to pairwise comparisons with
540 the central and southern US, which showed moderate
541 levels of genetic variation (Table 8), as well as the distant
542 placement of China 3 on the NJ tree (Fig. 6).

543 Of particular interest in the investigation of source
544 populations for the introduction of US kudzu is the NJ
545 tree clade composed of all the US subpopulations and
546 Japan 2, the centrally located Japanese subpopulation
547 (Fig. 6). With a bootstrap value of 50, these four subpop-
548 ulations can be distinguished from the rest of the tree
549 and within this clade, Japan 2 and US 3, the southern
550 US, are paired together with a support of 92. These find-
551 ings suggest that central Japan is a source of US kudzu.
552 Its association with US 3, the southern US populations,
553 makes sense considering that this area was where kudzu
554 was first planted for soil erosion control and where
555 farmers cultivated kudzu for fodder at the behest of the
556 US government. Our study is the first to provide mo-
557 lecular evidence to support the hypothesis of Japan as a
558 genetic source of US kudzu. However, a wider sampling
559 across the native Asian range coupled with higher num-
560 bers of genetic markers would increase statistical power

and confidence for testing genetic associations between 561
introduced and native kudzu, efforts that are currently 562
underway. 563

564 Conclusions

565 This study produced critical genomic resources for the
566 highly invasive kudzu vine by characterizing transcrip-
567 tomes and producing marker databases for SNP and
568 eSSR markers, foundational resources for understanding
569 ecological adaptation that may enable future insights
570 into invasiveness through gene discovery, marker-trait
571 analyses, and further genetic diversity studies. We exem-
572 plified the utility of our marker databases by assessing
573 the genetic diversity of native and introduced popula-
574 tions of kudzu using seven eSSRs. As a naturalized inva-
575 sive vine that was intentionally introduced throughout
576 millions of acres of the southeastern US, kudzu presents
577 unique challenges for management, especially given its
578 high genetic diversity across the US, a finding supported
579 by our genetic diversity analyses. The origin of this gen-
580 etic diversity remains a matter of speculation, however,
581 this study has begun to refine the proposed hypothesis
582 of single or multiple introductions from different genetic
583 populations. This study is the first to provide molecular
584 evidence that indicates the island of Honshu, Japan as
585 one source of US kudzu. Our analyses suggest either a
586 single introduction from a highly diverse source popula-
587 tion in Japan, or more likely multiple introductions from
588 multiple sources, potentially also from northern Japan
589 (Island of Hokkaido) or northern China. Given the eco-
590 logical and economic devastation wrought by kudzu in
591 the United States, it is critical that we improve our un-
592 derstanding of the history, process, origin(s), and im-
593 pacts of the U.S. kudzu invasion. We have assembled
594 transcriptomes and mined them for eSSRs that we have
595 provided as a resource for further genetic studies into
596 the origin(s) and range expansions of kudzu to that end.
597 By increasing both the sample ranges and sizes it should
598 be possible to identify more accurately the origin of
599 introduction and the number of introductions with the
600 markers we have developed, efforts that are currently
601 underway.

602 Methods

603 Plant material for transcriptome sequencing and 604 population genetics

605 Transcriptomic work in this study incorporated plant
606 tissues from two accessions of kudzu, *P. montana* var.
607 *lobata*, and one accession of tropical kudzu, *N. phaseo-*
608 *loides* [formerly *Pueraria phaseoloides* (Roxb.) Benth.].
609 One kudzu accession (noted here as Pmnk6) was wild
610 collected from Williamsburg, Virginia [voucher spec-
611 imen G. Tate s.n. (WILLI) collected 8 July 2013]. Leaf tis-
612 sue was collected in RNALater and preserved at -20 ° F

613 prior to RNA extraction. The other two plants were
614 grown from seed obtained from the United States De-
615 partment of Agriculture (USDA) Germplasm Resources
616 Information Network seed bank: accession PI 434246 of
617 *P. montana* var. *lobata* (noted here as CPP27) was field
618 collected in 1979 from the United States, locality un-
619 known, and is maintained by the Coffeerville Plant Mate-
620 rials Center, Soil Conservation Service, Coffeerville, MS;
621 accession PI 470272 of *N. phaseoloides* (noted here as
622 CPP02) was donated in 1981 from a field collection by
623 D.R. Bienz, 5 Jun 1981, Banjarbaru, S. Kalimantan,
624 Indonesia. Seeds were grown to maturity in the green-
625 house at Cornell University (Ithaca, NY, US) for 3 years
626 prior to RNA extraction. For eSSR screening and popu-
627 lation genetic studies, we sampled 75 accessions repre-
628 senting all three varieties of *P. montana* throughout
629 their native and US introduced range: US (25), China
630 (21), Japan (22) and Thailand (7) (Additional file 14).
631 Leaf material was immediately stored in silica for desic-
632 cation. Genomic DNA was extracted from samples using
633 Autogen robotics (Autogen Inc.) and a modified CTAB
634 extraction protocol [48].

635 RNA extraction and transcriptome sequencing

636 For the two accessions raised in the greenhouse, tissues
637 were flash frozen in liquid nitrogen prior to RNA extrac-
638 tion. *Neustanthus phaseoloides* (CPP02) was sampled for
639 young leaves, young shoot tips, and buds. Unfortunately,
640 kudzu never flowered in the greenhouse, so only young
641 shoot tips and young leaves were harvested for CPP27.
642 For the wild collected kudzu (Pmnk6), only young leaves
643 were harvested. RNA extraction, cDNA library construc-
644 tion, and transcriptome sequencing were carried out as
645 previously described [31]. cDNA libraries from CPP27
646 and CPP02 were multiplexed with two other libraries
647 not reported here across one titer plate on the Roche
648 454 Genome Sequencer FLX platform using Titanium
649 chemistry at the Brigham Young University Sequencing
650 Center (Provo, UT, US). Pmnk6 was also multiplexed
651 with three other transcriptomes not reported here and
652 sequenced using Roche 454 pyrosequencing, but using
653 Roche's next improvement on the titanium chemistry
654 that produced reads ~800 bp long. The raw sequence
655 data generated from CPP27, Pmnk6, and CPP02 were
656 deposited at the National Center for Biotechnology Infor-
657 mation (NCBI) Sequence Read Archive (SRA) under acces-
658 sion numbers SRR5925648, SRR5925647, and SRR5925649,
659 respectively.

660 De novo transcriptome assemblies

661 Raw reads were assessed for quality with FastQC [49]
662 and subsequently cleaned with ConDeTri [50], a
663 content-dependent read trimmer under the following
664 settings: reads below 50 bp were removed, Phred high

quality score thresholds (hq) were set to 25 and low
quality score thresholds (lq) were set to 10; the fraction of
bases per read having to exceed hq were set to 0.8 and the
minimum number of high quality bases (mh) and max-
imum number of low quality bases (ml) within the sliding
window were set to 30 and 5, respectively. Cleaned reads
were de novo assembled using Trinity (v2.0.6) [51] under
default parameters on two high-performance computing
clusters: the Smithsonian Institution High Performance
Cluster (SI/HPC) and the George Washington University
Colonial One Cluster. In order to minimize redundant
transcripts, a by-product of the assembly process,
CD-HIT-EST was used with a threshold of 0.9 to obtain
unique transcripts [52]. To evaluate the quality of the as-
semblies, criteria including the number of aligned reads,
total number of contigs produced, mean contig length,
N50, and transcript annotations were considered. RSEM
[53] and Bowtie2 [54] were used to identify the number of
aligned reads in the assembled transcriptomes. The KRA-
KEN suite was utilized in conjunction with prokaryote and
fungal databases to identify potential contaminants within
the transcriptomes [55]. BUSCO (v1.1b1), a pipeline used
to accurately annotate core genes in eukaryotic genomes,
was used to determine the completeness of the assemblies
[56]. At the time of use, BUSCO utilized a plant core data-
base of 956 single copy genes that are shared between *Ara-
bidopsis*, *Oryza*, *Populus*, and *Vitis* [57]. Reciprocal Best
BLAST Hits (RBH) between transcripts and among trans-
cripts were performed on a local installation of Galaxy
[58–60] and Toolshed [61] to characterize the number of
shared, homologous transcripts recovered in each Trinity
assembled transcriptomes [62, 63].

665 Functional annotation of transcriptomes

666 We used transcripts (contigs + singletons) assembled by
667 Trinity to annotate our transcriptomes (CPP27, CPP02,
668 and Pmnk6). To identify candidate coding regions, we
669 filtered sequences based on a minimum amino acid
670 length of 100 using the TransDecoder program v2.0.1
671 [64] with the TransDecoder.LongOrfs command. BlastP
672 and Pfam searches were carried out to detect open
673 reading frames (ORFs) with similarity to known pro-
674 teins and to maximize sensitivity for capturing ORFs
675 that may have functional significance. The BlastP
676 search was done using the Swissprot database with
677 the E-value of 1E-5 and Pfam search was done using
678 HMMER [65] and the Pfam database [66]. Output
679 files from the BlastP and Pfam searches were used to
680 ensure that peptides with BLAST or domain hits were
681 retained by running the TransDecoder.Predict com-
682 mand. The peptide sequences from the final candidate
683 ORFs were used to run BlastP searches against the
684 NCBI's nonredundant (nr) database with the E-value
685 of 1E-5 on the SI/HPC. The BLAST results were then
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717

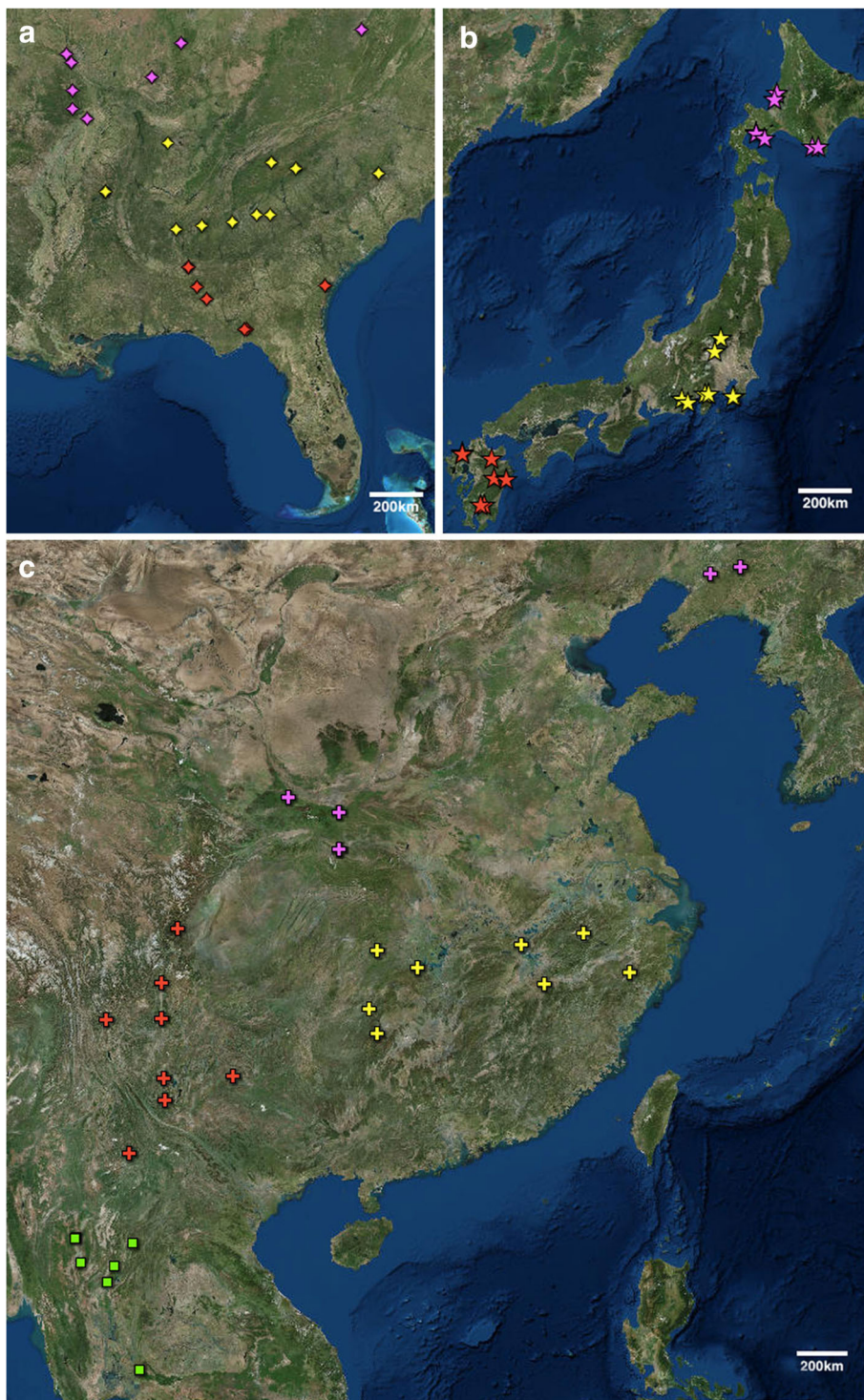


Fig. 7 Sampling sites: (a) United States: US 1, US 2, US 3 (25); (b) Japan: JP 1, JP2, JP 3 (22); and (c) China: CN 1, CN 2, CN 3 (21) and Thailand: TH (7)

f7.1
f7.2
f7.3

718 imported into the Blast2GO program v1.9.3 [67] to as-
719 sign Gene Ontology (GO) terms. We ran mapping, anno-
720 tation and InterProScan analyses for the three
721 transcriptomes separately.

Single nucleotide polymorphism identification

For SNP identification among the kudzu accessions, we used the transcripts (contigs + singletons) from our CPP27, Pmnnk6, and CPP02 assemblies and also incorporated two

722
723
724
725

publicly available *P. montana* var. *lobata* transcriptomes, SRX480408 from China derived from two tissues [68], and DRA001736 from Japan consisting of five pooled tissues [69]. We assembled the public sequences using Trinity as described above. Multiple pairwise comparisons between transcriptomes were conducted to evaluate the distribution of SNPs between US kudzu samples (CPP27 vs. Pmnk6) and identify intergeneric SNPs between kudzu and *N. phaseoloides* (CPP27 vs. CPP02 and Pmnk6 vs. CPP02). Additionally, the two US kudzu samples were combined by concatenating the two transcript files such that the samples represent the diversity in 'US kudzu' and subsequently compared to tropical kudzu to further identify intergeneric SNPs (CPP27/Pmnk6 vs. CPP02). Lastly, SNPs were called via comparison of all four *P. montana* var. *lobata* transcriptomes, with the transcriptome from Japan used as reference (Japan vs. CPP27/Pmnk6/China). The Japan transcriptome was chosen as reference because it incorporated the highest number of tissues, thus putatively having the higher chance of capturing greater expressed sequence diversity within the genome. To call SNPs, GS Reference Mapper v2.9 (454 Life Sciences, Roche, US) was used under default settings. The transcriptome composed of the greatest number of tissues was used as the reference to which reads from the others were assembled against. We used only high-confidence variants (454HCDiffs, >95%) in each comparison and further filtered these variants to those having 20× or greater coverage. To ensure the highest SNP call quality, we discarded any SNPs where 1) the reference or variant involved one or more N's or 2) the reference or variant allele was a single nucleotide insertion or deletion or did not include a point mutation in the length variant [70].

759 Expressed simple sequence repeat (eSSR) loci discovery, 760 screening and characterization

761 The ConDeTri cleaned, Trinity assembled, and
762 redundancy-vetted transcripts of CPP27 and CPP02
763 were mined for di-, tri-, tetra-, penta-, and hexanucleo-
764 tide microsatellites with MSATCOMMANDER [71].
765 Afterwards, MSATCOMMANDER and Primer3 [72]
766 were used to design primer pairs for each species with
767 an expected product size ranging from 100 to 450 bp.
768 Primer lengths were allowed to range from 18 to 22 bp,
769 annealing temperatures were optimized at 60 °C, and
770 GC contents were held between 30 and 70%. Developed
771 primers for both species were then cross-compared to
772 identify homologous primer regions, which could signify
773 interspecies transferability. The corresponding tran-
774 scripts for primers that were shared between *P. lobata*
775 and *N. phaseoloides* were blasted against the GenBank
776 nonredundant database using BLASTX [73] with an
777 *E*-value of 10^{-10} to determine the function of their associ-
778 ated unigenes. Pmnk6, SRX480408 [68] and DRA001736

[69] transcriptomes were not utilized for eSSR discovery 779
because none were available at the time eSSR mining took 780
place. Thirty potential eSSR primer pairs were chosen 781
from those discovered here and initially screened against a 782
subset of accessions (Additional file 11). Seventeen of the 783
30 primer pairs represent putatively homologous eSSRs 784
present in both *P. montana* var. *lobata* and *N. phaseo-* 785
loides (primer pairs designated PP) while the rest are *P.* 786
montana var. *lobata* specific (primer pairs designated PL). 787
The method of Culley et al. [19] was used to screen, 788
optimize and amplify eSSRs. Primer pairs were eliminated 789
based on the Culley et al. [19] protocol if they produced 790
superfluous primer diming between the specific and tailed 791
primers or produced PCR products of unexpected size. 792
Primer pairs were further eliminated if 1) primers did not 793
amplify viable product as seen via gel electrophoresis, 2) 794
primers amplified more bands than expected, or 3) 795
primers were monomorphic. 796

797 Screening of primer pairs against a subset of seven ac- 797
cessions ultimately yielded seven primer pairs that were 798
characterized across all 75 accessions. Primers, fluores- 799
cent dyes, and Culley method tail adaptors used for each 800
of the seven eSSRs are listed in Table 5. Initial rounds of 801
amplification across the entire sampling set were per- 802
formed in 12 µL reactions containing 1X Biolase NH₄ 803
buffer, 1.0 µL primer mix, 1.2 mM MgCl₂, 0.12 µL of 804
8 µM dNTPs, 0.35 U of Taq polymerase (Biolase), and 805
5-80 ng DNA template. PCR was performed on an Ap- 806
plied Biosystems 2720 thermocycler with settings of 95C 807
for 5 min, followed by 35 cycles of 95C for 30s, 50C for 808
45 s, 72C for 30s, and a final 72C extension for 5 min. 809
Annealing temperatures were adjusted between 810
51.5C-58C for primers PP13, PL1, PL11, and PP2. Prod- 811
uct bands were resolved using 1.5% sodium borate gels 812
containing GelRed stain and visualized under UV light. 813
Accessions that failed to amplify after two or more initial 814
attempts were subsequently attempted with an adjusted 815
concentration of 2.38 µg MgCl₂ per reaction. Further 816
failed amplifications were then tried using AmpliTaq Gold 817
using reaction mix 1X AmpliTaq buffer, 1.0µL of primer 818
mix, 2.86 µg MgCl₂, 1.2µL of 8 µM dNTPs, 0.375 U of 819
AmpliTaq Gold Taq polymerase [0.075 µL of 1000 U in 820
200 µL], and 5-80 ng DNA template. Successful products 821
were genotyped using an ABI3730 sequencer at the 822
Smithsonian NMNH LAB facilities. Genotypes were called 823
using GeneMapper (v5.0) [74]. 824

825 Examination of population structure and genetic diversity 826 indices

827 Genetic population structuring was assessed with
828 STRUCTURE v2.3.4 [75] and STRUCTURE HAR-
829 VESTER v0.6.94 [76]. The length of burnin period was
830 set to 100,000, while the number of MCMC reps after
831 burnin was set to 900,000, resulting in a total of 1

832 million generations. No LOCPRIOR information was
 833 provided for the STRUCTURE runs. A job consisting of
 834 10 iterations, evaluating Ks from 1 to 10 for the 75 *P.*
 835 *montana* accessions, was run and the results were
 836 uploaded to STRUCTURE Harvester for analysis. The
 837 optimal K was assessed via the Evanno et al. [77]
 838 method. Individual and population files were loaded into
 839 CLUMPP v1.1.2 [78] to address label switching and the
 840 potential for multimodality across the 10 STRUCTURE
 841 iterations. The CLUMPP program utilized the FullSearch
 842 method, the number of individuals in each population
 843 influenced weights, and the pairwise matrix similarity
 844 statistic was set to G'. All additional options remained as
 845 default settings. CLUMPP outputs for the individual and
 846 population files were visualized with DISTRUCT v1.1
 847 [79]. Genetic diversity statistics were calculated in Arle-
 848 quin v3.5.1.9 [80]. The default parameters of Arlequin
 849 were used on our 75-individual data set that was subdivi-
 850 ded from the four sampled nations to 10 geographic-
 851 ally defined subpopulations: US (3), China (3), Japan (3),
 F7 852 and Thailand (1) (Fig. 7). The subpopulation designa-
 853 tions were based primarily on geographic proximity that
 854 allowed for groupings of at-least five individuals along
 855 similar latitudinal lines; however, due to the different
 856 scales of sampling done across nations, the ranges of the
 857 latitudinal boundaries of the subpopulations differed.
 858 POPTREEW [81] was used to make a neighbor joining
 859 (NJ) distance tree with F_{st} distances [82] for the above
 860 listed subpopulations. Bootstrap support for the tree was
 861 calculated with 10,000 replicates.

862 Additional files

863
 864 **Additional file 1: Table S1.** Trinity contig reads mapped back to the
 865 raw and cleaned reads. Numbers of cleaned and raw reads mapped back
 866 to contigs via Bowtie2. (PDF 126 kb)
 867
 868 **Additional file 2: Table S2.** Contaminated reads as assessed by Kraken.
 869 Number (percentage) of cleaned reads annotated by Kraken as
 870 prokaryotic or fungal. (PDF 126 kb)
 871
 872 **Additional file 3: Figure S1.** CPP27 Top-Hit Species Distribution. Top-hit
 873 species distribution of CPP27 proteins annotated against NCBI's non-
 874 redundant database showing the highest distribution of hits against leg-
 875 ume species. (PDF 808 kb)
 876
 877 **Additional file 4: Figure S2.** Pmnk6 Top-Hit Species Distribution. Top-hit
 878 species distribution of Pmnk6 proteins annotated against NCBI's non-
 879 redundant database showing the highest distribution of hits against leg-
 880 ume species. (PDF 753 kb)
 881
 882 **Additional file 5: Figure S3.** CPP02 Top-Hit Species Distribution. Top-hit
 883 species distribution of CPP02 proteins annotated against NCBI's non-
 884 redundant database showing the highest distribution of hits against leg-
 885 ume species. (PDF 2607 kb)
 886
 887 **Additional file 6:** SNPs_Pmnk6_vs_CPP27. High-confidence single nu-
 888 cleotide polymorphisms between US kudzu accessions Pmnk6 (variant:
 Var) and CPP27 (reference: Ref). Accno: contig in reference; Pos: position;
 Nuc: nucleotide; Total Depth: number of variant reads aligned against the
 reference; Var Freq: frequency of variant SNP within aligned reads; # Fwd:
 number of forward reads with variant; # Rev.: number of reverse reads

with variant; # Fwd Total: number of forward-aligned reads total; # Rev.
 Total: number of reverse-aligned reads total. (XLSX 578 kb) 889
 890
Additional file 7: SNPs_CPP02_vs_CPP27. High-confidence single nu- 891
 cleotide polymorphisms between tropical kudzu CPP02 (reference: Ref) 892
 and kudzu accession CPP27 (variant: Var). Abbreviations as described for 893
 Additional file 6. (XLSX 4932 kb) 894
Additional file 8: SNPs_CPP02_vs_Pmnk6. High-confidence single nu- 895
 cleotide polymorphisms between tropical kudzu CPP02 (reference: Ref) 896
 and kudzu accession Pmnk6 (variant: Var). Abbreviations as described for 897
 Additional file 6. (XLSX 11073 kb) 898
Additional file 9: SNPs_CPP02_vs_Pmnk6_CPP27. High-confidence single 899
 nucleotide polymorphisms between tropical kudzu CPP02 (reference: 900
 Ref) and a composite transcriptome comprising reads from kudzu acces- 901
 sions CPP27 and Pmnk6 (variant: Var). Abbreviations as described for Add- 902
 itional file 6. (XLSX 11520 kb) 903
Additional file 10: SNPs_Japan_vs_Pmnk6_CPP27_China. High- 904
 confidence single nucleotide polymorphisms among kudzu accessions 905
 from Japan (reference: Ref) and reads from US kudzu (Pmnk6 and CPP27) 906
 and China (variants: Var). Abbreviations as described for Additional file 6. 907
 (XLSX 30817 kb) 908
Additional file 11: Table S3. Thirty primer pairs tested for polymorphic 909
 amplification in *Pueraria montana*. Primers labeled PP were designed 910
 from kudzu and tropical kudzu transcriptomes whereas those designated 911
 PL were designed from kudzu only. Bold primers are those used for 912
 population genetic analyses in this study. F: forward primer; R: reverse 913
 primer. (PDF 33 kb) 914
Additional file 12: Figure S4. Delta K of STRUCTURE run (K = 3). Plot of 915
 Delta K for STRUCTURE analyses from K = 2 through K = 9, with K = 3 seen 916
 as the optimal number of genetic clusters. (PDF 18 kb) 917
Additional file 13: Table S4. Allele table for *Pueraria* subpopulations. 918
 Number of alleles discovered for each locus within each subpopulation, 919
 with mean and standard deviation (SD) for each subpopulation and each 920
 locus. (PDF 19 kb) 921
Additional file 14: Table S5. Plant material used for eSSR validation 922
 and population genetics. Species determination, subpopulation 923
 designation (pop), country and state/province/island of origin within the 924
 United States (US), China (CN), Japan (JP) or Thailand (TH), voucher 925
 information, accession number, and geographical coordinates for each of 926
 the 75 plants used in the population genetic analyses. (PDF 34 kb) 927
 928

Abbreviations

BLAST: Basic local alignment search tool; bp: Base pair; 929
 BUSCO: Benchmarking universal single-copy orthologs; eSSR: Expressed 930
 simple sequence repeat; GO: Gene ontology; hq: High quality; lq: Low 931
 quality; mh: Minimum high quality; ml: Maximum low quality; NCBI: National 932
 Center for Biotechnology Information; nr: Nonredundant; ORF: Open reading 933
 frame; RBH: Reciprocal best hits; RIN: RNA integrity; SI/HPC: Smithsonian 934
 Institution High Performance Cluster; SNP: Single nucleotide polymorphism; 935
 SRA: Sequence read archive; SSR: Simple sequence repeat; TNR: Tri- 936
 nucleotide repeat 937
 938

Acknowledgements

We thank Susan Sherman-Broyles and Jane L. Doyle for help in sustaining 939
 plants in the greenhouse and to Beth Chambers and Gus Tate, Herbarium of 940
 the College of William and Mary, for help in obtaining voucher specimens 941
 for Pmnk6. Additionally, we thank Cheng-Xin Fu, Lu-Xian Liu, Xin-fen Gao 942
 and Bo Xu for assistance collecting in China, Tetsukazu Yahara, Tadashi Kajita, 943
 Firouzeh Javadi, Tomoko Otao and Yumi Kagawa for help in Japan, and Vora- 944
 dol Chamchumroon, Kongkanda Chayamarit, Thaveechok Jumruschay Rum- 945
 rada Meeboonya, Nannapat Pattharahirantricit, Rachun Pooma, Sukontip 946
 Sirmongkol, and Ruth P. Clark for help in Thailand. Computations were com- 947
 pleted in part on the Smithsonian Institution High Performance Cluster (SI/ 948
 HPC) and the George Washington University Colonial One Cluster. We also 949
 thank the Computational Biology Institute at the George Washington Univer- 950
 sity for graduate support for MSH. 951
 952

953 Funding

954 This research was supported by funding from the US National Science
955 Foundation to ANE (DEB-1352217) and JJD (DEB-0948800).

956 Availability of data and materials

957 The transcriptomes generated and analyzed during the current study are
958 available in the NCBI repository, [Study PRJNA397892, accessions:
[Q3] 959 SRR5925647, SRR5925648, and SRR5925649, [http://www.ncbi.nlm.nih.gov/
960 bioproject/397892](http://www.ncbi.nlm.nih.gov/bioproject/397892), release date 30 June 2018]. The SNP data generated
961 during this study are included in this published article's Additional files 6, 7,
962 8, 9 and 10 however, the SNPs contained in Additional file 10 are not
963 publicly available due to file size restrictions but are available from the
964 corresponding author upon reasonable request.

965 Authors' contributions

966 All authors contributed to various aspects of this work (ordered by degree of
967 contribution): conceived the study (ANE, MSH); aided in experimental design
968 (MSH, ANE); obtained research funds (ANE, JJD); coordinated activities (ANE,
969 MSH); obtained and grew plants (ANE, MSH, JJD); RNA Isolation and Library
970 Prep (ANE); transcriptome assembly and analyses (MSH, MV, ANE);
971 microsatellite primer design (MSH); microsatellite primer validation (MSH, GM,
972 DZ, RZMR); prepared figures (MSH, MV, ANE); contributed to preparation of
973 the manuscript (MSH, ANE, MV, JJD, KAC). All authors edited and approved
974 the final manuscript.

975 Ethics approval and consent to participate

976 All plant material was collected in accordance with institutional, national,
977 and international guidelines and under appropriate permits. Permits and
978 voucher specimens are deposited at the US National Herbarium (US) with all
979 specimens determined by Dr. Ashley N. Egan.

980 Competing interests

981 The authors declare that they have no competing interests.

982 Publisher's Note

983 Springer Nature remains neutral with regard to jurisdictional claims in
984 published maps and institutional affiliations.

985 Author details

[Q2] 986 ¹Department of Biology, George Washington University, Washington, DC,
987 USA. ²Computational Biology Institute, Milken Institute School of Public
988 Health, George Washington University, Washington, DC, USA. ³Department of
989 Botany, National Museum of Natural History, Smithsonian Institution,
990 Washington, DC, USA. ⁴Present address: College of Engineering, Oregon
991 State University, Corvallis, OR, USA. ⁵Present address: Department of Biology,
992 Washington University in St. Louis, St. Louis, MO, USA. ⁶Present address:
993 Department of Biology, Indiana University Bloomington, Bloomington, IN,
994 USA. ⁷School of Integrated Plant Science, Plant Breeding and Genetics
995 Section, Cornell University, Ithaca, NY, USA. ⁸Department of Invertebrate
996 Zoology, National Museum of Natural History, Smithsonian Institution,
997 Washington, DC, USA.

998 **Received: 25 April 2018 Accepted: 15 May 2018**

999 

1000 References

- 1001 1. Pueraria V d MLJG. Botanical characteristics. In: Keung WM, editor. *Pueraria:*
1002 *the genus Pueraria*. New York: Taylor and Francis; 2002. p. 1–28.
- 1003 2. Egan AN, Pan B. Resolution of polyphyly in *Pueraria* (Leguminosae,
1004 Papilionoideae): the creation of two new genera, *Haymondia* and
1005 *Toxicopueraria*, the resurrection of *Neustanthus*, and a new combination in
1006 *Teyleria*. *Phytotaxa*. 2015;218:201–26.
- 1007 3. Van der Maesen LJG. Revision of the genus *Pueraria* with some notes on
[Q4] 1008 *Teyleria* (Leguminosae): Taylor & Francis; 1985.
- 1009 4. Van der Maesen LJG. *Pueraria*, the kudzu and its relatives, an update of the
1010 taxonomy. In: Sorensen M, editor. Proceedings of the first international
1011 symposium on tuberous legumes, Gualdeloupe, FWI. Denmark: DSR
1012 Boghandel; 1994. p. 55–86.
- 1013 5. Egan AN, Vatanparast M, Cagle W. Parsing polyphyletic *Pueraria*: delimiting
1014 distinct evolutionary lineages through phylogeny. *Mol. Phylogenet. Evol.*
1015 2016;104:44–59.
6. Forseth IN Jr, Innis AF. Kudzu (*Pueraria montana*): history, physiology, and 1016
ecology combine to make a major ecosystem threat. *Crit Rev Plant Sci.* 1017
2004;23:401–13. 1018
7. Follak S. Potential distribution and environmental threat of *Pueraria lobata*. 1019
Cent Eur J of Biol. 2011;6:457–69. 1020
8. Westbrooks R. Invasive plants, changing the landscape of America: fact 1021
book. In: Federal Interagency Committee for the Management of Noxious 1022
and Exotic Weeds: Washington; 1998. 1023
9. Kudzu SD. In: Simberloff D, Rejmanek D, editors. Encyclopedia of biological 1024
invasions. California: University of California Press; 2011. p. 396–9. 1025
10. Pappert RA, Hamrick JL, Donovan LA. Genetic variation in *Pueraria lobata* 1026
(Fabaceae), an introduced, clonal, invasive plant of the southeastern 1027
United States. *Am J Bot.* 2000;87:1240–5. 1028
11. Dk J, Jiang CJ, Britton KO, Sun JH, Tang J. Characterizing specimens of 1029
kudzu and related taxa with RAPDs. *Castanea.* 2003;68:254–60. 1030
12. Sun JH, Li Z-C, Jewett DK, Britton KO, Ye WH, Ge X-J. Genetic diversity of 1031
Pueraria lobata (kudzu) and closely related taxa as revealed by inter-simple 1032
sequence repeat analysis. *Weed Res.* 2005;45:255–60. 1033
13. Heider B, Fischer E, Berndt T, Schultze-Kraft R. Analysis of genetic variation 1034
among accessions of *Pueraria montana* (Lour.) Merr. *Var. lobata* and *Pueraria*
1035 *phaseoloides* (Roxb) Benth. Based on RAPD markers. *Genet Resour Crop*
1036 *Evol.* 2007;54:529–42. 1037
14. Hoffberg SL, Bentley KE, Lee JB, Myhre KE, Iwao K, Glenn TC, et al. 1038
Characterization of 15 microsatellite loci in kudzu (*Pueraria montana var.*
1039 *lobata*) from the native and introduced ranges. *Conserv Genet Resour.* 2015;7:
1040 403–5. 1041
15. Bentley K, Mauricio R. High degree of clonal reproduction and lack of large-scale 1042
geographic patterning mark the introduced range of the invasive vine, kudzu
1043 (*Pueraria montana var. lobata*) in North America. *Am J Bot.* 2016;103:1499–507. 1044
16. Strickler SR, Bombarely A, Mueller LA. Designing a transcriptome next- 1045
generation sequencing project for a nonmodel plant species. *Am J Bot.*
1046 2012;99:257–66. 1047
17. Varshney RK, Sigmund Rm Borner A, Korzun V, Stein N, Sorrells ME, et al. 1048
Interspecific transferability and comparative mapping of barley EST-SSR
1049 markers in wheat, rye, and rice. *Plant Sci.* 2005;168:195–202. 1050
18. Ellis J, Burke J. EST-SSRs as a resource for population genetic analyses. 1051
Heredity. 2007;99:125–32. 1052
19. Culley TM, Stamper TI, Stokes RL, Brzyski JR, Hardiman NA, Klooster MR, et al. 1053
An efficient technique for primer development and application that integrates
1054 fluorescent labeling and multiplex PCR. *Appl Plant Sci.* 2013;1:1–10. 1055
20. Wright S. Evolution and the genetics of populations. Vol. 4. Variability within 1056
and among natural populations. Chicago: University of Chicago Press; 1978. 1057
21. Stewart CN Jr, Tranel PJ, Horvath DP, Anderson JV, Rieseberg LH, Westwood 1058
JH, et al. Evolution of weediness and invasiveness: charting the course for
1059 weed genomics. *Weed Sci.* 2009;57:451–62. 1060
22. Ekblom R, Galindo J. Applications of next-generation sequencing in 1061
molecular ecology of non-model organisms. *Heredity.* 2011;107:1–15. 1062
23. Egan AN, Schlueter J, Spooner DM. Applications of next-generation 1063
sequencing in plant biology. *Am J Bot.* 2012;99:175–85. 1064
24. Wen J, Egan AN, Dikow RB, Zimmer EA. Utility of transcriptome sequencing 1065
for phylogenetic inference and character evolution. In: Hörandl E,
1066 Appelhans MS, editors. Next-generation sequencing in plant systematics.
1067 Königstein: Koeltz scientific books; 2015. p. 51–91. 1068
25. Ilut DC, Coate JE, Luciano AK, Owens TG, May GD, Farmer A, et al. A 1069
comparative transcriptomic study of an allotetraploid and its diploid
1070 progenitors illustrates the unique advantages and challenges of RNA-Seq in
1071 plant species. *Am J Bot.* 2012;99:383–96. 1072
26. Grover CE, Salmon A, Wendel JF. Targeted sequence capture as a powerful 1073
tool for evolutionary analysis. *Am J Bot.* 2012;99:312–9. 1074
27. Egan AN, Doyle J. A comparison of global, gene-specific, and relaxed clock 1075
methods in a comparative genomics framework: dating the polyploid
1076 history of soybean (*Glycine max*). *Syst Biol.* 2010;59:534–47. 1077
28. Cannon SB, McKain MR, Harkess A, Nelson MN, Dash S, Deyholos MK, et al. 1078
Multiple polyploid events in the early radiation of nodulating and
1079 nonnodulating legumes. *Mol. Biol. Evol.* 2015;32:193–210. 1080
29. Glover NM, Redestig H, Dessimoz C. Homoeologs: what are they and how 1081
do we infer them? *Trends Plant Sci.* 2016;21:609–21. 1082
30. Champigny MJ, Sung WW, Catana V, Salwan R, Summers PS, Dudley SA, 1083
et al. RNA-Seq effectively monitors gene expression in *Eutrema salsugineum*
1084 plants growing in an extreme natural habitat and in controlled growth
1085 cabinet conditions. *BMC Genomics.* 2013;14:578. 1086

- 1087 31. Vatanparast M, Shetty P, Chopra SP, Doyle JJ, Sathyanarayana N, Egan AN. 1088 Transcriptome sequencing and marker development in winged bean 1089 (*Psophocarpus tetragonolobus*; Leguminosae). *Sci Rep.* 2016;6:29070. 1158
- 1090 32. Kudapa H, Azam S, Sharpe AG, Taran B, Li R, Deonovic B, et al. 1091 Comprehensive transcriptome assembly of chickpea (*Cicer arietinum* L.) 1092 using sanger and next generation sequencing platforms: development and 1093 applications. *PLoS One.* 2014;9:e86039. 1162
- 1094 33. Sudheesh S, Sawbridge TI, Cogan NO, Kennedy P, Forster JW, Kaur S. De 1095 novo assembly and characterization of the field pea transcriptome using 1096 RNA-Seq. *BMC Genomics.* 2015;16:611. 1163
- 1097 34. Khalturin K, Hemmrich G, Fraune S, Augustin R, Bosch TCG. More than just 1098 orphans: are taxonomically-restricted genes important in evolution? *Trends 1099 Genet.* 2009;25:404–13. 1164
- 1100 35. Kumar S, Banks TW, Cloutier S. SNP discovery through next-generation 1101 sequencing and its applications. *Int J Plant Genomics.* 2012; <https://doi.org/10.1155/2012/831460>. 1165
- 1102 36. Varshney RK, Chabane K, Hendre PS, Aggarwal RK, Graner A. Comparative 1103 assessment of EST-SSR, EST-SNP and AFLP markers for evaluation of genetic 1104 diversity and conservation of genetic resources using wild, cultivated and elite 1105 barleys. *Plant Sci.* 2007;173:638–49. 1166
- 1106 37. Zalapa JE, Cuevas H, Zhu H, Steffan S, Senalik D, Zeldin E, et al. Using next- 1107 generation sequencing approaches to isolate simple sequence repeat (SSR) 1108 loci in the plant sciences. *Am J Bot.* 2012;99:193–208. 1167
- 1109 38. Morgante M, Hanafey M, Powell W. Microsatellites are preferentially associated 1110 with nonrepetitive DNA in plant genomes. *Nat Genet.* 2002;30:194–200. 1168
- 1111 39. Lopez L, Barreiro R, Fischer M, Koch MA. Mining microsatellite markers from 1112 public expressed sequence tags databases for the study of threatened 1113 plants. *BMC Genomics.* 2015;16:781. 1169
- 1114 40. Thumilan BM, Sajeevan RS, Biradar J, Madhuri T, Nataraja KN, Sreeman SM. 1115 Development and characterization of genic SSR markers from Indian 1116 mulberry transcriptome and their transferability to related species of 1117 Moraceae. *PLoS One.* 2016;11:e0162909. 1170
- 1118 41. Wang P, Yang L, Zhang E, Qin Z, Wang H, Liao Y, et al. Characterization and 1119 development of EST-SSR markers from a cold-stressed transcriptome of 1120 Centipedegrass by Illumina paired-end sequencing. *Plant Mol Biol Rep.* 1121 2017;35:215–23. 1171
- 1122 42. Yang Z, Peng ZS, Yang H. Identification of novel and useful EST-SSR markers 1123 from de novo transcriptome sequence of wheat (*Triticum aestivum* L.). 1124 *Genet Mol Res.* 2016;15 1172
- 1125 43. Kartzinel TR, Hamrick JL, Wang C, Bowsher AW, Quigley BG. Heterogeneity of 1126 clonal patterns among patches of kudzu, *Pueraria montana* var. *lobata*, an 1127 invasive plant. *Ann Bot.* 2015;116:739–50. 1173
- 1128 44. Ellstrand NC, Roose ML. Patterns of genotypic diversity in clonal plant 1129 species. *Am J Bot.* 1987;74:123–31. 1174
- 1130 45. Balloux F, Lehmann L, de Meeùs T. The population genetics of clonal and 1131 partially clonal diploids. *Genetics.* 2003;164:1635–44. 1175
- 1132 46. Halkett FJ, Simon JC, Balloux FO. Tackling the population genetics of clonal 1133 and partially clonal organisms. *Trends Ecol. Evol.* 2005;20:194–201. 1176
- 1134 47. Shurtleff W, Aoyagi A. The book of kudzu: a culinary and healing guide. 1135 Brookline: Autumn Press; 1997. 1177
- 1136 48. Doyle JJ, Doyle JL. A rapid DNA isolation procedure for small quantities of 1137 fresh leaf tissue. *Phytochem Bull.* 1987;19:11–5. 1178
- 1138 49. Andrews S. FastQC. A quality control tool for high throughput sequence data. 1139 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>. Accessed 14 1140 Aug 2017. 1179
- 1141 50. Smeds L, Künstner A. ConDeTri – a content dependent read trimmer for 1142 Illumina data. *PLoS One.* 2011;6:e26314. 1180
- 1143 51. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Full- 1144 length transcriptome assembly from RNA-Seq data without a reference 1145 genome. *Nat Biotechnol.* 2011;29:644–52. 1181
- 1146 52. Fu L, Niu B, Zhu Z, Wu S, Li W. CD-HIT: accelerated for clustering the next- 1147 generation sequencing data. *Bioinformatics.* 2012;28:3150–2. 1182
- 1148 53. Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data 1149 with or without a reference genome. *BMC Bioinformatics.* 2011;12:323. 1183
- 1150 54. Langmead B, Salzberg SL. Fast gapped-read alignment with bowtie 2. *Nat 1151 Methods.* 2012;9:357–9. 1184
- 1152 55. Davis MPA, van Dongen S, Abreu-Goodger C, Bartonicek N, Enright AJ. Kraken: 1153 a set of tools for quality control and analysis of high-throughput sequence 1154 data. *Methods.* 2013;63:41–9. 1185
- 1155 56. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 1156 BUSCO: assessing genome assembly and annotation completeness with 1157 single-copy orthologs. *Bioinformatics.* 2015;31:3210–2. 1158
- 1159 57. Duarte JM, Wall PK, Edger PP, Landherr LL, Ma H, Pires PK, et al. 1160 Identification of shared single copy nuclear genes in *Arabidopsis*, *Populus*, 1161 *Vitis*, and *Oryza* and their phylogenetic utility across various taxonomic 1162 levels. *BMC Evol Biol.* 2010;10:61. 1163
- 1163 58. Goecks J, Nekrutenko A, Taylor J. The galaxy team. Galaxy: a comprehensive 1164 approach for supporting accessible, reproducible, and transparent 1165 computational research in the life sciences. *Genome Biol.* 2010;11:R86. 1166
- 1166 59. Blankenberg D, Von Kuster G, Coraor N, Ananda G, Lazarus R, Mangan M, 1167 et al. Galaxy: a web-based genome analysis tool for experimentalists. *Curr 1168 Protoc Mol Biol.* 2010; 1169
- 1169 60. Giardine B, Riemer C, Hardison RC, Burhans R, Elmtski L, Shah P, et al. Galaxy: 1170 a platform for interactive large-scale genome analysis. *Genome Res.* 2005;15: 1171
- 1171 61. Blankenberg D, Von Kuster G, Bouvier E, Baker D, Afgan E, Stoler N, et al. 1172 Dissemination of scientific software with galaxy toolshed. *Genome Biol.* 1173 2014;15:403. 1174
- 1174 62. Cock PJA, Chilton JM, Grüning B, Johnson JE, Soranzo N. NCBI BLAST+ 1175 integrated into galaxy. *GigaScience.* 2015;4:39. 1176
- 1176 63. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos, Bealer K, et al. 1177 BLAST+: architecture and applications. *BMC Bioinformatics.* 2009;10:421. 1178
- 1178 64. Haas B, Papanicolaou A. TransDecoder <https://transdecoder.github.io>. 1179 Accessed 14 Aug 2017. 1179
- 1179 65. Finn RD, Clements J, Eddy SR. HMMER web server: interactive sequence 1180 similarity searching. *Nucleic Acids Res.* 2011;39:W29–37. 1181
- 1181 66. Sonnhammer ELL, Eddy SR, Durbin R. Pfam: a comprehensive database of 1182 protein domain families based on seed alignments. *Proteins.* 1997;28:405–20. 1183
- 1182 67. Conesa A, Götz S, García-Gómez J, Terol J, Talón M, Robles M. Blast2GO: a 1183 universal tool for annotation, visualization and analysis in functional 1184 genomics research. *Bioinformatics.* 2005;21:3674–6. 1185
- 1184 68. Wang X, Li S, Li J, Li C, Zhang Y. De novo transcriptome sequencing in 1185 *Pueraria lobata* to identify putative genes involved in isoflavones 1186 biosynthesis. *Plant Cell Rep.* 2015;34:733–43. 1187
- 1186 69. Han R, Takahashi H, Nakamura M, Yoshimoto N, Suzuki H, Shibata D, et al. 1187 Transcriptomic landscape of *Pueraria lobata* demonstrates potential for 1188 phytochemical study. *Front Plant Sci.* 2015;6:426. 1189
- 1188 70. Schmutz J, McClean PE, Mamidi S, Wu GA, Cannon SB, Grimwood J, et al. A 1189 reference genome for common bean and genome-wide analysis of dual 1190 domestications. *Nat Genet.* 2014;46:707–13. 1191
- 1190 71. Faircloth BC. MSATCOMMANDER: detection of microsatellite repeat arrays and 1191 automated, locus-specific primer design. *Mol Ecol Resour.* 2008;8:92–4. 1192
- 1191 72. Rozen S, Skaletsky H. Primer3 on the WWW for general users and for 1192 biologist programmers. *Methods Mol Biol.* 2000;132:365–86. 1193
- 1192 73. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. 1193 Gapped BLAST and PSI-BLAST: a new generation of protein database search 1194 programs. *Nucleic Acids Res.* 1997;25:3389–402. 1195
- 1193 74. Chatterji S, Pachter L. Reference based annotation with GeneMapper. 1194 *Genome Biol.* 2006;7:R29. 1196
- 1194 75. Pritchard JK, Stephens M, Donnelly P. Inference of population structure 1195 using multilocus genotype data. *Genetics.* 2000;155:945–59. 1197
- 1195 76. Earl DA, vonHoldt BM. STRUCTURE HARVESTER: a website and program for 1196 visualizing STRUCTURE output and implementing the Evanno method. 1197 *Conserv Genet Resour.* 2012;4:359–61. 1198
- 1196 77. Evanno G, Regnaut S, Goudet J. Detecting the number of clusters of 1197 individuals using the software STRUCTURE: a simulation study. *Mol Ecol.* 1198 2005;14:2611–20. 1199
- 1197 78. Jakobsson M, Rosenberg NA. CLUMPP: a cluster matching and permutation 1199 program for dealing with label switching and multimodality in analysis of 1200 population structure. *Bioinformatics.* 2007;23:1801–6. 1201
- 1198 79. Rosenberg NA. DISTRUCT: a program for the graphical display of population 1202 structure. *Mol Ecol Notes.* 2004;4:137–8. 1203
- 1199 80. Excoffier L, Laval G, Schneider S. Arlequin (version 3.0): an integrated 1204 software package for population genetics data analysis. *Evol Bioinforma.* 1205 2005;1:47–50. 1206
- 1200 81. Takezaki N, Nei M, Tamura K. POPTREE: Web version of POPTREE for 1207 constructing population trees from allele frequency data and computing 1208 other population statistics. *Mol Biol Evol.* 2014;31:1622–4. 1209
- 1201 82. Latter BDH. Selection in finite populations with multiple alleles. III. Genetic 1210 divergence with centripetal selection and mutation. *Genetics.* 1972;70:475–90. 1211
- 1202 1227

Author Query Form

Journal: BMC Genomics

Title: De novo transcriptome assembly of *Pueraria montana* var. *lobata* and *Neustanthus phaseoloides* for the development of eSSR and SNP markers: narrowing the US origin(s) of the invasive kudzu

[Q1]

Authors: Matthew S. Haynsen, Mohammad Vatanparast, Gouri Mahadwar, Dennis Zhu, Roy Z. Moger-Reischer, Jeff J. Doyle, Keith A. Crandall, Ashley N. Egan

Article: 4798

Dear Authors,

During production of your paper, the following queries arose. Please respond to these by annotating your proofs with the necessary changes/additions. If you intend to annotate your proof electronically, please refer to the E-annotation guidelines. We recommend that you provide additional clarification of answers to queries by entering your answers on the query sheet, in addition to the text mark-up.

Query No.	Query	Remark
Q1	Author names: Please confirm that the author names are presented accurately and in the correct sequence (given names/initials, family name). Author 1: Given name: Matthew Given name: S. Family name: Haynsen Author 2: Given name: Mohammad Family name: Vatanparast Author 3: Given name: Gouri Family name: Mahadwar Author 4: Given name: Dennis Family name: Zhu Author 5: Given name: Roy Given name: Z. Family name: Moger-Reischer Author 6: Given name: Jeff Given name: J. Family name: Doyle Author 7: Given name: Keith Given name: A. Family name: Crandall Author 8: Given name: Ashley	

Query No.	Query	Remark
	Given name: N. Family name: Egan	
Q2	Please check if the affiliations are presented correctly.	
Q3	URL: Please check that the following URLs are working. If not, please provide alternatives: http://www.ncbi.nlm.nih.gov/bioproject/397892	
Q4	Citation details for Reference [3] is incomplete. Please supply the "publisher location" of this reference. Otherwise, kindly advise us on how to proceed.	
Q5	Citation details for Reference [42] is incomplete. Please supply the "page range" of this reference. Otherwise, kindly advise us on how to proceed.	
Q6	URL: Please check that the following URLs are working. If not, please provide alternatives: http://www.bioinformatics.babraham.ac.uk/projects/fastqc/	
Q7	Citation details for Reference [59] is incomplete. Please supply the "volume number and page range" of this reference. Otherwise, kindly advise us on how to proceed.	
Q8	URL: Please check that the following URLs are working. If not, please provide alternatives: https://transdecoder.github.io	