

1 **Title: Transcribed microsatellites often influence gene expression in natural sunflower**
2 **populations**

3 Authors: Chathurani Ranathunge^{1*}, Gregory L. Wheeler^{1,2}, Melody E. Chimahusky¹, Andy D.
4 Perkins³, Sreepriya Pramod^{1,4} and Mark. E. Welch¹

5 ¹Department of Biological Sciences, Mississippi State University, Starkville MS 39762.

6 ²Current Address: Department of Evolution, Ecology, and Organismal Biology, 1735 Neil
7 Avenue, Ohio State University, Columbus, OH 43210

8 ³Department of Computer Science and Engineering, Mississippi State University, Starkville MS
9 39762

10 ⁴Current Address: Altria Client Services LLC, 601 E Jackson Street, Richmond VA 23219

11

12 * Authors to whom all correspondence should be addressed: car494@msstate.edu

13 Running title: Microsatellites influence gene expression in sunflower

14 Keywords: *Helianthus annuus*, gene expression, microsatellite, sunflower

15

16

17

18

19

20

21 **ABSTRACT**

22 Microsatellites are common in most species. While an adaptive role for these highly
23 mutable regions has been considered, little is known concerning their contribution towards
24 phenotypic variation. Here we used populations of the common sunflower (*Helianthus annuus*)
25 at two latitudes to quantify the effect of microsatellite allele length on phenotype at the level of
26 gene expression. We conducted a common garden experiment with seed collected from
27 sunflower populations in Kansas and Oklahoma followed by an RNA-Seq experiment on 95
28 individuals. The effect of microsatellite allele length on gene expression was assessed across
29 3325 microsatellites that could be consistently scored. Our study revealed 479 microsatellites at
30 which allele length significantly correlates with gene expression (significant expression short
31 tandem repeats or eSTRs). Further, when irregular allele sizes not conforming to the motif length
32 were removed from the analysis, the number of eSTRs rose to 2379. The percentage of variation
33 in gene expression explained by eSTRs ranged from 1 – 86% when controlling for population
34 and allele-by-population interaction effects at the 479 eSTRs initially identified. Further, 70.4%
35 of these eSTRs are located in untranslated regions (UTRs). A Gene Ontology (GO) analysis
36 revealed that eSTRs are significantly enriched for GO terms associated with cis- and trans-
37 regulatory processes. These findings suggest that a substantial number of transcribed
38 microsatellites can influence gene expression levels. This result is consistent with the hypothesis
39 that these repetitive elements can serve as engines of adaptation.

40

41

42

43 INTRODUCTION

44 The molecular basis of adaptation is of fundamental concern to evolutionary geneticists.
45 Adaptation can occur via selection acting on standing genetic variation, or by that acting on
46 novel mutations. While standing genetic variation might allow populations to respond rapidly to
47 novel selective pressures in the short-term, novel beneficial mutations may arise at a finite rate.
48 The relative contributions of these two processes have been argued (Orr 2010). What has been
49 particularly puzzling to many theorists is the abundance of genetic variation in natural
50 populations despite the expectation that its loss due to selection and drift should eliminate it.
51 Ultimately, limited availability of heritable variation should serve as a form of evolutionary
52 constraint. In spite of these factors that should limit the availability of additive genetic variation,
53 there is much evidence to suggest that numerous traits can continue to respond to selection even
54 after many generations of intense directional selection (Dudley and Lambert 1992; Mackay
55 1995; Yoo et al. 1980). These empirical results are consistent with the existence of mechanisms
56 that continuously generate additive genetic variation at rates commensurate with its loss through
57 selection.

58 These observations led Barton (1990) to suggest the existence of an abundant source of
59 non-deleterious mutations that could affect quantitative traits. Kashi et al (1997) further
60 emphasized that to be considered advantageous mutators with significant contributions toward
61 rapid adaption, these mutations should be widely dispersed across genomes, associated with
62 functional regions as regulatory elements or function as components of coding regions, and
63 undergo high mutation rates leading to quantitative effects on phenotypes. These suggested
64 favorable features are consistent with those observed in highly mutable microsatellites.

65 Microsatellites, or simple tandem repeats (STRs), are genomic regions consisting of short
66 motifs that are 1- 6 bp in length, tandemly repeated up to a few dozen times (Vogt 1990).
67 Microsatellites show high indel mutation rates (10^{-2} - 10^{-6} / generation) (Ellegren 2004) resulting
68 from mechanisms that may include replication slippage (Tautz and Renz 1984). Their mutation
69 rates are estimated to be several orders of magnitude greater than that observed in non-repetitive
70 DNA (Jarne and Lagoda 1996; Li 1997). Despite their apparent fit for the role of advantageous
71 mutators, microsatellites have long been perceived as neutrally evolving, non-functional regions,
72 and have been used as the “molecular marker of choice” in population genetics and forensics
73 (Jarne and Lagoda 1996). This long-standing perception was questioned by the abundance of
74 highly mutable microsatellites in structural regions placing them in positions that could influence
75 gene function and gene products (Li et al. 2002; Li et al. 2004). Further, non-random distribution
76 of microsatellites in genomes has been reported in several organisms, including fruit fly
77 (*Drosophila melanogaster*) (Bachtrog et al. 1999), thale cress (*Arabidopsis thaliana*), rice (*Oryza*
78 *sativa*) (Lawson and Zhang 2006; Morgante et al. 2002), and common sunflower (*Helianthus*
79 *annuus* L.) (Pramod et al. 2014). Microsatellite variation in functional regions has been
80 consequently linked to phenotypic variation. Notable examples include microsatellites linked to
81 human neurodegenerative diseases such as fragile X syndrome (Verkerk et al. 1991) and
82 Huntington’s disease (Andrew et al. 1993).

83 The claim that microsatellites can generate adaptive genetic variation is now supported
84 by a growing list of studies. Research has linked microsatellites to variation in skeletal
85 morphology in domesticated dogs (Fondon and Garner 2004), social behavior changes in voles
86 (Hammock and Young 2005) and some primates (Hopkins et al. 2012), pathogenesis in bacteria
87 (Moxon et al. 1994; Moxon et al. 2006), plasticity in adherence to substrates in *Saccharomyces*

88 *cerevisiae* (Verstrepen et al. 2005), and thermal sensitivity in *Drosophila melanogaster* (Costa et
89 al. 1991) among others. However, little is known about the mechanisms by which microsatellites
90 can influence phenotypes. An intriguing model proposed to explain the underlying mechanisms
91 by which microsatellites may influence phenotypes is the “tuning knob” model. This model
92 likens microsatellites to tuning knobs where stepwise changes in microsatellite allele lengths can
93 have stepwise effects on phenotypes, allowing selection to adjust or fine tune a population’s
94 phenotypes corresponding to environmental stresses (Kashi and King 2006; King et al. 1997;
95 Trifonov 2004). The model explains that these effects can be mediated either by modulating gene
96 expression or by facilitating structural changes in proteins.

97 Microsatellites located upstream of genes as well as those in 5’ untranslated regions
98 (UTRs) and introns might affect gene expression while those in 3’ UTRs are more likely to
99 influence transcript stability (Li et al. 2004). Microsatellites located in coding regions are likely
100 to generate structural changes in proteins, but may also play a regulatory role as well (Gemayel
101 et al. 2010; Li et al. 2004). The role of microsatellites in modulating gene expression has been
102 supported empirically by introducing microsatellites of different lengths to promoter regions of
103 genes in *Saccharomyces cerevisiae* (Lee and Maheshri 2012; Vinces et al. 2009) and in Tausch's
104 goatgrass (*Aegilops tauschii*) (Ryan et al. 2010). Experimentally introduced microsatellite tracts
105 in coding regions of genes have also been utilized to show their role in altering protein structures
106 in *Arabidopsis thaliana* (Golubov et al. 2010). It is apparent that many studies have focused on
107 the effect of microsatellites on specific genes, as opposed to investigating the global role that
108 microsatellites may play across genomes potentially leading to adaptive evolution. Such large-
109 scale genome level studies focusing on potential phenotypic effects of multiple microsatellite
110 loci across populations have been few and sporadic (Fahima et al. 2002; Gymrek et al. 2015;

111 Nevo et al. 2005). Here, we attempt to test the relevance of the tuning knob model at the genome
112 level using transcribed sequences. Microsatellites located within transcribed regions are
113 particularly accessible in this regard given that RNA-Seq data can be used to both genotype
114 transcribed microsatellites, and allow for estimation of allele-specific expression at
115 microsatellite-encoding transcripts.

116 In this study we use natural populations of the common sunflower (*Helianthus annuus* L.)
117 transecting a latitudinal cline from Kansas to Oklahoma. We chose sunflower as a model system
118 given their adaptability to diverse environmental conditions across their broad geographical
119 range in North America (Heiser et al. 1969). Particularly, populations across latitudes
120 demonstrate heritable clinal variation in a number of traits including flowering time (Blackman
121 et al. 2011) and seed oil content (Linder 2000). These distinct patterns of variation in adaptive
122 traits across latitudinal populations provide an opportunity to test the hypothesis that
123 microsatellites generate some of the adaptive genetic variation found within and among natural
124 populations.

125 From a probabilistic viewpoint, it is reasonable to expect that transcribed microsatellites,
126 due to their proximity to functional regions, are likely to influence gene expression as *cis*-
127 regulatory elements. We utilized an RNA-Seq approach to assess the degree to which allele
128 lengths of transcribed microsatellites influence gene expression across multiple transcripts to
129 identify potential tuning knobs. We predicted that allele lengths of transcribed microsatellites
130 that function as tuning knobs are correlated with gene expression levels.

131

132

133 **RESULTS**

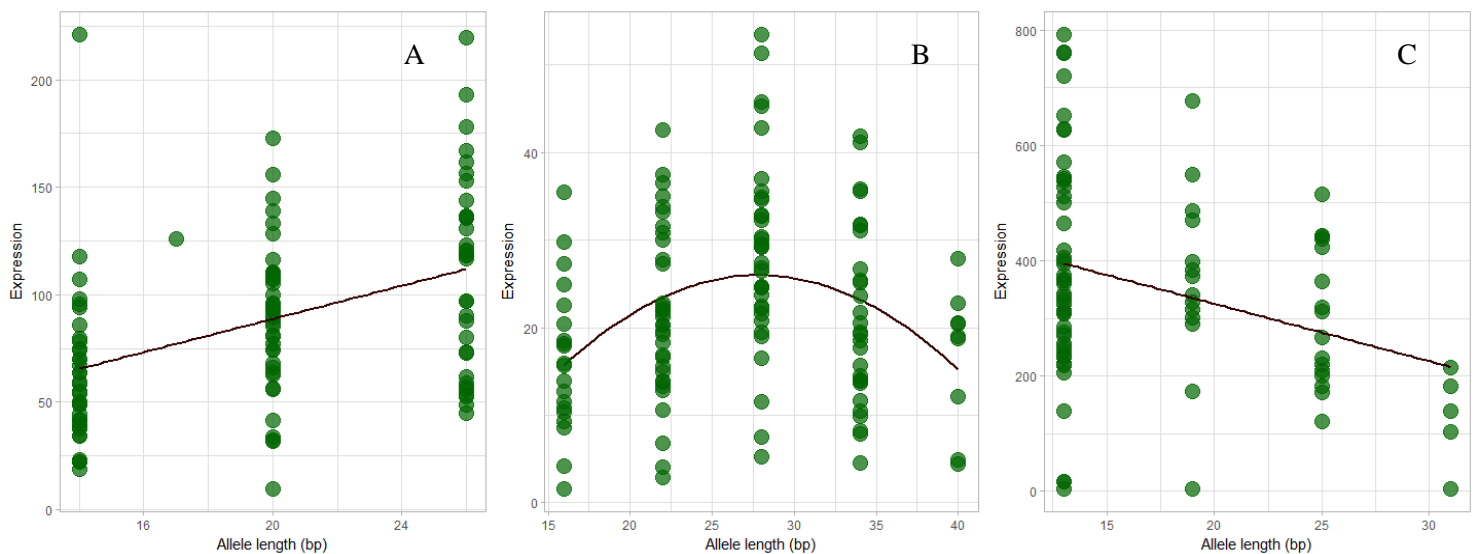
134 *Microsatellite search*

135 Tandem Repeat Finder identified 19,104 potential microsatellites in the reference
136 transcriptome. We were able to genotype 3,325 of these microsatellites in 2,640 transcripts
137 across consistently using RepeatSeq. A total of 685 (25.9%) transcripts harbored more than one
138 microsatellite. Hexanucleotides (38.89%) were identified as the most abundant microsatellite
139 motif length, followed by trinucleotides (29.17%) (Supplemental Table S6; Supplemental
140 Fig.S1). A total of 355 standardized motif types were identified. Of the 355 different motif types
141 identified, the ACC repeat motif has the highest frequency (6.32%), followed by ATC repeats
142 (5.78%) (Supplemental Table S7; Supplemental Fig.S2).

143 *Microsatellites with significant allele length effect on gene expression (eSTRs)*

144 When individuals with irregular allele sizes were included in the ANCOVA, 816 (24.5%)
145 of the microsatellites scored showed a significant allele length effect on gene expression (log
146 transformed) ($p < 0.05$), controlling for population and allele-by-population interaction effects.
147 After correcting for multiple comparisons (Storey 2002), a total number of 479 (14.4%)
148 microsatellites in 449 unique transcripts were identified as eSTRs (Supplemental Table S8). The
149 percent of variation in gene expression explained by allele lengths ranged from 1% - 86% in the
150 479 eSTRs. ANCOVA results for 237 microsatellites suggest a positive relationship between
151 allele length and gene expression levels while 242 microsatellites exhibited a negative
152 correlation between allele length and gene expression (Figure 1). When the irregular allele sizes
153 that did not conform to the motif length were removed from the analysis, the number of eSTRs
154 rose to 2379 (71.5% of the genotyped microsatellites) after correcting for multiple comparisons
155 (Storey 2002). The effect of allele length on gene expression ranged from 0.077 to 79.9% in the

156 2379 microsatellites (Supplemental Table S9). While we acknowledge that the first approach
157 could have been conservative in identifying eSTRs, removal of all irregular allele sizes from the
158 analysis could have resulted in excessive manipulation of data as demonstrated by the inflated
159 estimates of correlation between allele length and gene expression. Taking these concerns into
160 consideration, we used the 479 eSTRs identified with the first approach for downstream
161 analyses.

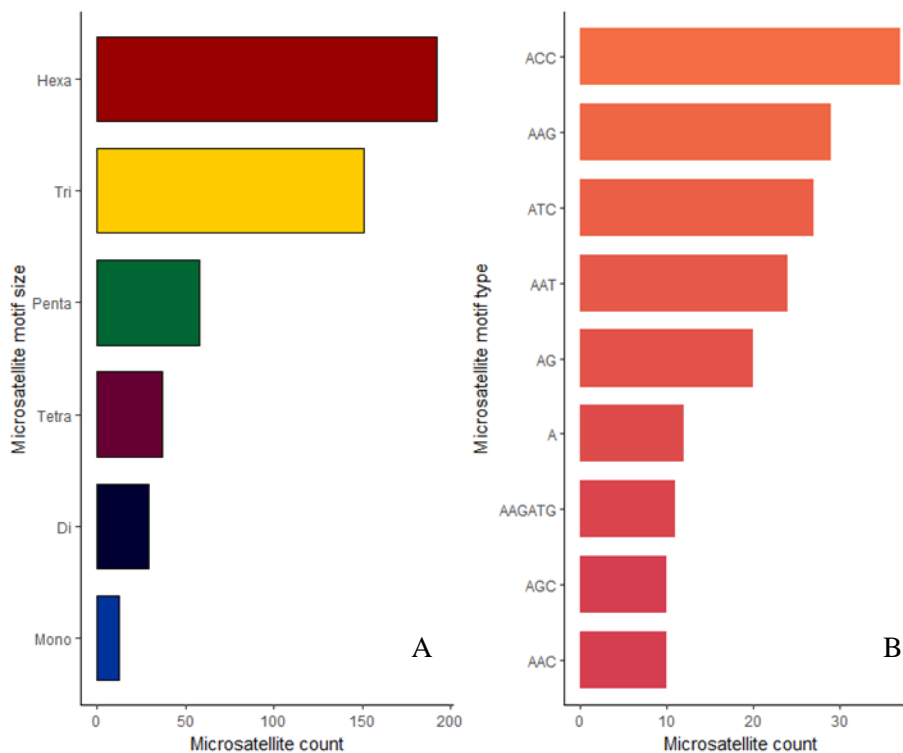


162 **Figure 1. The effect of microsatellite allele length on gene expression.** An Analysis of
163 Covariance (ANCOVA) revealed different patterns of correlation between microsatellite allele
164 length and allele specific gene expression (eSTR). (A), (B) and (C) show positive, quadratic and
165 negative correlation patterns observed between allele length (bp) and allele specific gene
166 expression (read counts per 100 million reads) in microsatellite loci located in contigs,
167 comp26672 (CCTTCT), comp49389 (GAACCA) and comp25013 (GACGGT), respectively.

168

169 We characterized the motif lengths associated with eSTR-containing transcripts in
170 addition to the location of the microsatellite relative to start and stop codons. Hexanucleotides
171 (40.1%) are the most abundant motif length within the 479 eSTRs followed by trinucleotides
172 (31.5%) (Figure 2; Supplemental Table S10). A total of 169 different motifs (standardized) were
173 identified within 449 transcripts (0.38 motif/transcript) containing the 479 eSTRs. ACC (21.9%)

174 was the most abundant motif within eSTRs followed by AAG (17.2%) (Figure 2; Supplemental
175 Table S11). Prior work detected significant enrichment of A and AG motif-containing
176 microsatellites within genes that are differentially expressed between the two latitudes in Kansas
177 and Oklahoma (Ranathunge *et al.* 2018). This suggested that A and AG repeat motifs are likely
178 to be involved in gene expression divergence among these sunflower populations. To test
179 whether specific motif-containing microsatellites are more likely to function as tuning knobs for
180 gene expression, we estimated enrichment of different microsatellite motif types within eSTRs
181 compared to the remaining 2846 microsatellites genotyped (non-eSTRs). However, we did not
182 detect any significant enrichment of specific motif types within eSTRs compared to non-eSTRs
183 (Chi-squared test, $p > 0.05$).

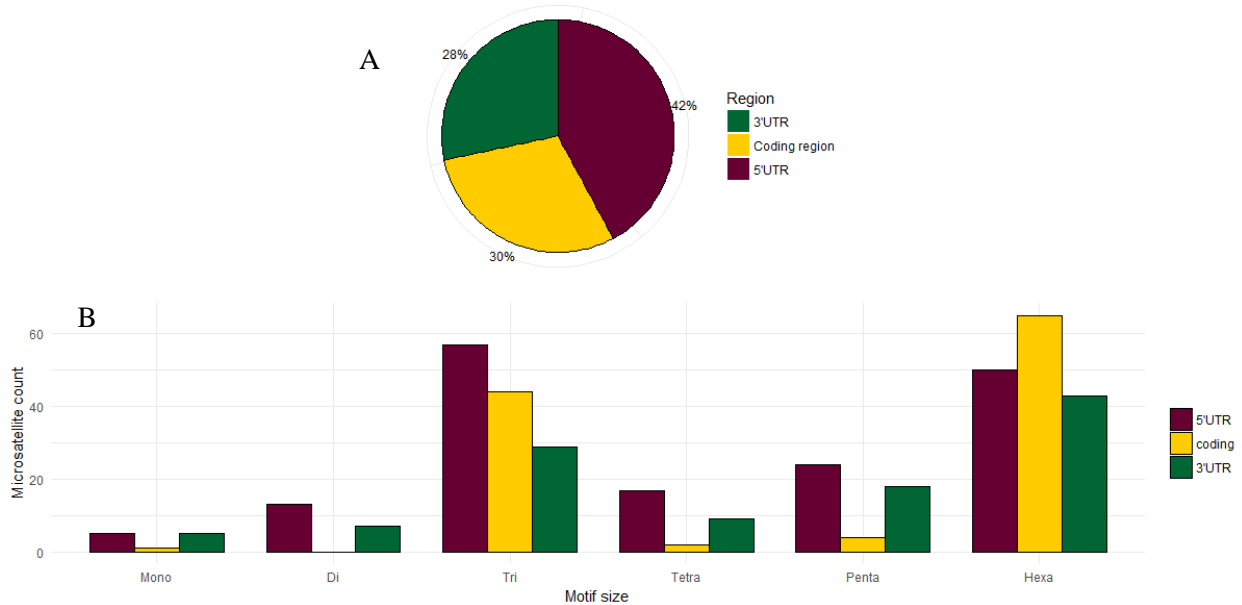


184

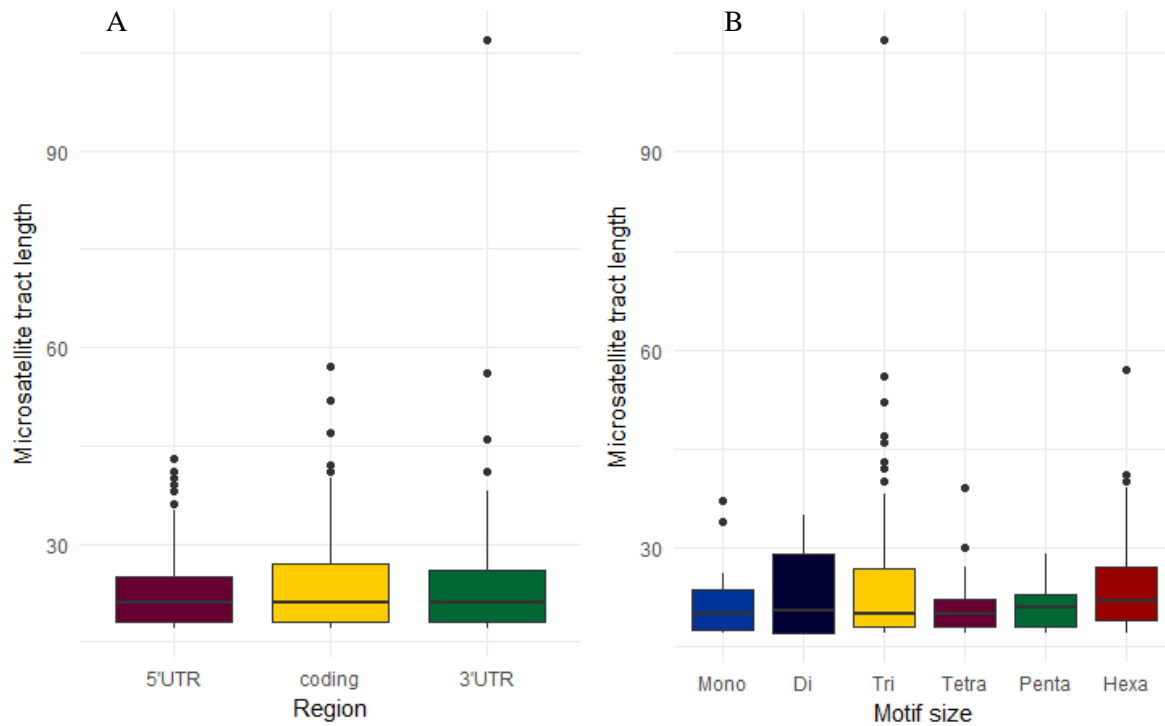
185 **Figure 2. The distribution of microsatellite motif sizes and types in eSTRs.** (A) Number of
186 mono-, di-, tri-, tetra-, penta- and hexanucleotides identified within eSTRs. (B) The top nine
187 microsatellite motif types and their counts within eSTRs.

188 eSTRs were most abundant within 5'UTRs (42.1%) followed by coding regions (29.6%)
189 and 3'UTRs (28.3%) (Figure 3). Of the 165 eSTRs located within the 5'UTRs, the most
190 abundant motif size present were trinucleotides (34.55%) followed by hexanucleotides (30.30%)
191 (Figure 3). Within the 116 eSTRs in the coding regions, hexanucleotides were the most abundant
192 motif size (56.03%) followed by trinucleotides (37.93%) while dinucleotides were absent in the
193 region. Mononucleotides (0.86%), tetranucleotides (1.92%), and pentanucleotides (0.96%) were
194 also scarce within coding regions (Figure 3). Hexanucleotides were also the most abundant
195 within the 111 eSTRs located within the 3'UTRs (38.74%) followed by trinucleotides (26.13%)
196 (Figure 3). ACC was the most abundant motif in all three regions; 5'UTR (8.48%), coding region
197 (8.62%) and 3'UTR (8.11%) (Supplemental Table S12). Prior work revealed that a transcript
198 containing a microsatellite within the 3'UTR is more likely to be differentially expressed
199 between the latitudinal populations in Kansas and Oklahoma (Ranathunge et al. 2018). In line
200 with this prediction, we tested whether a transcript containing a microsatellite within 5'UTR,
201 coding, or 3'UTR is more likely to function as tuning knobs for gene expression. We assessed
202 the enrichment of eSTRs within 5'UTRs, coding regions, and 3'UTRs in comparison to
203 frequency of non-eSTRs within the three regions. We did not detect any significant enrichment
204 of eSTRs within any of the three regions (Chi-squared test, $p > 0.05$). Further, there was no
205 significant difference in mean microsatellite tract lengths among eSTRs located in 5' UTRs,
206 coding regions and 3' UTRs (KW test, $p = 0.45$) (Figure 4). No significant differences were
207 detected in eSTR tract length among different motif sizes (KW test, $p = 0.07$) (Figure 4).

208



209 **Figure 3. The distribution of eSTRs located within the three regions; 5'UTR, coding and**
210 **3'UTR. (A) Pie chart represents the percentage of eSTRs located within the three regions;**
211 **5'UTR, coding and 3'UTR. (B) The counts of different eSTR motif sizes identified within each**
212 **region.**



213 **Figure 4. The variation in eSTR tract lengths. (A) The variation in eSTR tract length by**
214 **region. (B) The variation in eSTR tract length by motif size.**

215 ***Validation of gene expression estimates by Real Time PCR (qPCR)***

216 All seven qPCR assays for eSTR-containing transcripts showed strong log-linear
217 relationships between C_T values and cDNA concentrations. Coefficient of determination
218 estimates (R^2) ranged from 0.97 – 0.98 and efficiencies (b) ranged from 0.71 to 0.9
219 (Supplemental Table S4). Regression analyses conducted with log C: R ratios for pairs of high
220 copy number and low copy number loci revealed strong correlation between loci for relative
221 concentrations estimated with qPCR and RNA-Seq data with R^2 values ranging from 0.77 to 0.96
222 (Supplemental Table S13; Supplemental Fig.S3).

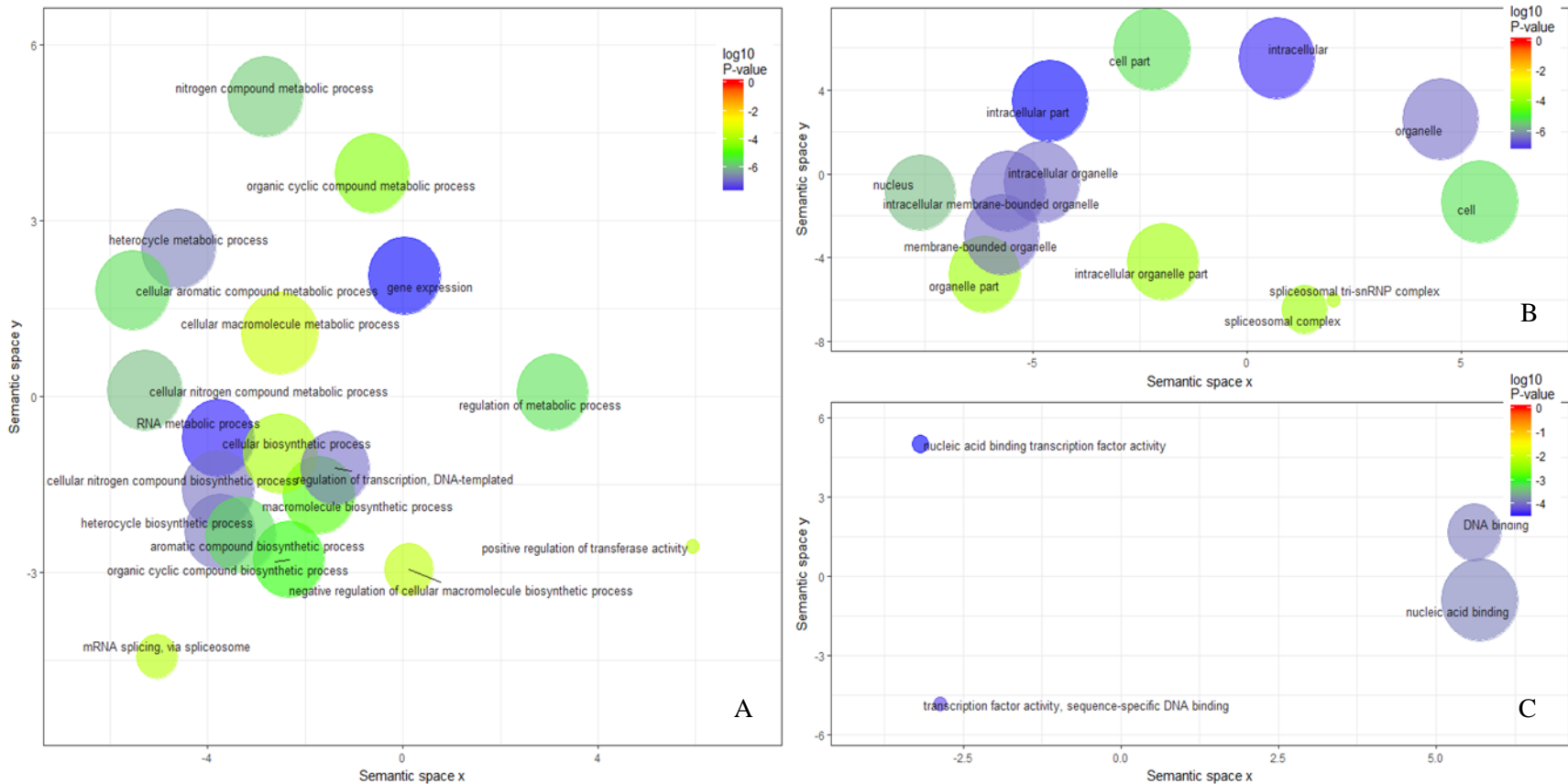
223 ***Validation of RNA-Seq derived microsatellite genotypes with fragment analysis***

224 The number of matches and mismatches in genotypes derived from fragment analysis and
225 RepeatSeq were counted to assess the concordance between the two methods for each locus.
226 Scoring was limited to 47, 77 and 92 individuals for the three loci, comp 47327 (A), comp 25013
227 (GACGGT) and comp 47993 (TTCAA) respectively based on genotype data availability for both
228 methods. The results indicate that the three microsatellite loci exhibited concordance between the
229 two genotyping methods. The three loci, comp47327, comp25013 and comp47993 showed 71%,
230 79% and 85% matches and 13%, 14% and 7% mismatches respectively. The missing data
231 percentages for the three loci were 16%, 7% and 8% respectively. The mismatch and missing
232 data percentages estimated for the three loci may represent errors arising from both RNA-Seq
233 and fragment analysis methods during variant calling.

234 ***Functional annotation and Gene Ontology (GO) enrichment analysis***

235 The BLASTX search against the *H. annuus* protein database identified unique hits for
236 17,985 contigs in the reference transcriptome (Supplemental Table S14). Of the 449 eSTR-

237 containing transcripts, 371 had blast hits that passed the significance threshold (Supplemental
238 Table S15). In comparison to GO terms associated with the annotated 17,985 genes in the
239 reference transcriptome, eSTR-containing transcripts were significantly enriched for 57 GO
240 terms (Supplemental Table 16). The simplified enriched GO term list included eight specific GO
241 terms (Table 1; Figure 5). They were classified under biological process (4), molecular function
242 (2) and cellular component (2). The most enriched GO term within the eSTR-containing
243 transcripts was associated with regulation of transcription (GO:0006355) in the biological
244 process category (Table 1; Figure 5). Within the molecular function category, the GO term
245 represented by most sequences in our list was transcription factor activity (GO:0003700) while
246 the GO term, spliceosomal complex (GO:005681) was the most overrepresented within the
247 cellular component category.



248 **Figure 5. Gene Ontology (GO) terms enriched within eSTRs.** GO enrichment analysis of eSTR-containing transcripts was
 249 conducted against a background of all expressed transcripts. The enriched terms were visualized using REVIGO. (A). Reduced GO
 250 terms related to biological processes. (B) Reduced GO terms associated with cellular components. (C) Reduced GO terms associated
 251 with molecular functions. The circle size represents the frequency of the GO term and the color represents the \log_{10} P-value based on
 252 the Fisher's exact test for enrichment.

253 **Table 1. Gene Ontology (GO) terms enriched within eSTR-containing transcripts in**
 254 ***Helianthus annuus* (reduced to most specific terms)**

GO ID	GO Name	GO Category	FDR	p - value	% in test set	% in background set
GO:0006355	regulation of transcription, DNA-templated	BP	6.93E-05	1.06E-07	17.73%	8.66%
GO:0003700	transcription factor activity, sequence-specific DNA binding	MF	0.01169	5.86E-05	10.47%	5.12%
GO:0003677	DNA binding	MF	0.024173	0.000129	15.12%	8.84%
GO:0005681	spliceosomal complex	CC	0.025453	0.000145	2.33%	0.41%
GO:0000398	mRNA splicing, via spliceosome	BP	0.027348	0.000158	2.62%	0.53%
GO:0051347	positive regulation of transferase activity	BP	0.034644	0.000204	1.45%	0.14%
GO:2000113	negative regulation of cellular macromolecule biosynthetic process	BP	0.035718	0.000215	6.69%	2.85%
GO:0097526	spliceosomal tri-snRNP complex	CC	0.037228	0.000232	1.16%	0.07%

255 Note: BP = Biological process, MF = Molecular function, CC = Cellular component

256

257 **DISCUSSION**

258 Here we tested specific predictions of the tuning knob model that proposes stepwise
259 effects of microsatellite allele length on phenotypes. This study focused on transcribed
260 microsatellites because both genotypic and phenotypic variation at the level of gene expression
261 can be assessed using RNA-Seq data. We used *Helianthus annuus*, common sunflower, seed
262 collected from two latitudes to infer the scope of this proposed mechanism for rapid evolutionary
263 change in natural populations. The seeds were germinated, and plants were grown in a common
264 garden experiment to minimize environmental variation. Further, we employed an experimental
265 design that included replicates from each of the two latitudes which allowed us to estimate the
266 relative effects of the microsatellite allele length on gene expression given the effects of other
267 components such as local genetic background and environmental variation within latitudes on
268 phenotypic variation. These populations also show latitudinal variation in a number of
269 morphological traits including flowering time, and plant height at budding that indicates the
270 heritable nature of phenotypic traits in these populations (Ranathunge et al. 2018). These
271 observations provide an impetus to study the underlying mechanisms of adaptation across this
272 latitudinal cline.

273 Using the RNA-Seq data from these populations, we were able to both estimate gene
274 expression levels and consistently genotype microsatellites in 2,640 transcripts containing 3325
275 microsatellites. Motif type and length searches across the 3325 microsatellites showed frequency
276 estimates consistent with similar surveys conducted on eukaryotic genomes (Qin et al. 2015;
277 Tóth et al. 2000) including those conducted on the *Helianthus annuus* transcriptome with the use
278 of an EST database (Pramod et al. 2014). These results demonstrate that the microsatellites used
279 to test the predictions of the tuning knob hypothesis in this study substantially represent the

280 composition of microsatellites in eukaryotic genomes (Supplemental Table 6; Supplemental
281 Table 7).

282 Our findings from the ANCOVAs estimating the effect of microsatellite allele length on
283 gene expression provide support for the tuning knob model at a large number of microsatellite
284 encoding transcripts. Of the loci scored for both gene expression level and microsatellite
285 genotypes (3325), 479 microsatellite loci (14.4%) (eSTRs) showed significant allele length effect
286 on gene expression. It is important to note that this first analysis was conducted without filtering
287 out irregular allele sizes that did not conform to the repeat motif. When irregular alleles were
288 removed from the analysis, the number of eSTRs rose to 2379 (71.5%). The irregular allele
289 lengths observed in some individuals could represent imperfect microsatellites resulting from
290 substitutions and indels within the microsatellite tract or sequencing artifacts. Typically,
291 imperfect microsatellites have been noted for lower mutation rates and higher stability in
292 comparison to perfect microsatellites (Kunkel and Soni, 1988). Sequence interruptions have been
293 previously reported in microsatellites linked to trinucleotide-repeat disorders (Chung *et al.* 1993;
294 Eichler *et al.* 1994; Kunst *et al.* 1997). Based on these findings, it is reasonable to expect that
295 mechanisms underlying functional microsatellites are likely to be affected by the level of
296 stability attributed to these repeat tracts. Perhaps irregularities in microsatellites can hinder the
297 cis-regulatory activity of these genetic elements.

298 Downstream analyses were limited to the 479 eSTRs identified with the more
299 conservative of the two approaches. The 479 eSTRs showed on average, 1% – 86% allele length
300 effect on gene expression when accounted for population and allele-by-population interaction
301 effects. The majority of the eSTRs identified in ours study showed a linear relationship between
302 microsatellite allele length and gene expression. Previous studies have also demonstrated

303 positive and negative linear relationships between microsatellite length and gene expression
304 (Contente et al. 2002; Gymrek et al. 2015; Shimajiri et al. 1999). We also tested whether our data
305 fit a quadratic model as previous experimental studies have demonstrated similar patterns of
306 correlation between experimentally constructed promoter microsatellite lengths and gene
307 expression in yeast (Vinces et al. 2009). Based on the significance threshold we set when
308 conducting the ANCOVA, 171 eSTRs showed support for a quadratic relationship between
309 microsatellite allele length and gene expression. However, when we examined the data for those
310 loci, most of these trends were identified as artifacts of rare allele lengths represented in the data
311 set. These findings suggest that the relationship between gene expression and extant alleles in
312 natural populations can typically be modeled as linear relationships.

313 The position of the microsatellite tract within genes may also provide insights regarding
314 potential mechanisms by which microsatellites modulate gene expression. To better understand
315 these different mechanisms, we estimated the frequencies of different motif types within the
316 eSTRs, and also examined the likely locations for the eSTRs within *H. annuus* genes. A previous
317 study conducted on the RNA-Seq data generated from the same populations indicated that
318 microsatellites with motif types A and AG as well as microsatellites located within the 3'UTRs
319 are likely to cause gene expression divergence among these latitudinal populations (Ranathunge
320 et al. 2018). In the current study we did not detect any significant enrichment of specific
321 microsatellite motif types within the eSTRs in comparison to non-eSTRs which suggests that the
322 motif type may not affect the microsatellite's ability to function a tuning knob. These contrasting
323 patterns observed between the study on microsatellites involved in differential gene expression
324 (Ranathunge et al. 2018) and the current study on eSTRs suggests the involvement of different
325 groups of microsatellites in generating large and subtle differences in gene expression.

326 Furthermore, except for one microsatellite locus in transcript comp29399, there were no
327 microsatellites in common between those identified as eSTRs in the current study and the
328 microsatellites identified within differentially expressed genes between these two latitudinal
329 populations by Ranathunge et al. 2018.

330 When we examined the location of the eSTRs within the transcribed region, our results
331 suggest that most eSTRs (42.1%) are located in 5'UTRs. Combined, microsatellites in 5' and 3'
332 UTRs accounted for 70.4% of the eSTRs. However, when we estimated the enrichment of eSTRs
333 in the three regions in comparison to the distribution of non-eSTRs genotyped, we did not detect
334 any significant difference. This suggests that the likelihood of microsatellites functioning as
335 tuning knobs may not depend on their location within the transcribed region. Several
336 experimental studies have demonstrated that some microsatellites in 5'UTRs are vital for the
337 expression of the gene (Kumar and Bhatia 2016; Streelman and Kocher 2002; Toutenhoofd et al.
338 1998). Some of the proposed mechanisms by which 5'UTR microsatellites may function as cis-
339 regulatory elements involve serving as transcription factor binding sites (Kumar and Bhatia
340 2016). Microsatellites in coding regions may also play a role in gene expression regulation. In
341 our study, 29.6% of the eSTRs were located in coding regions. Variation in coding region
342 microsatellites is linked to changes in the structure of proteins including transcription factors as
343 documented in several studies (Fondon and Garner 2004; Gemayel et al. 2015; Lee et al. 2003).
344 However, the mechanisms by which they may function as cis-regulatory elements directly
345 influencing gene expression levels are not entirely clear. Some triplet repeats including those in
346 coding regions are proposed to affect nucleosome binding thereby affecting transcription rates
347 (Sandman and Reeve 1999; Wang 2007). Coding region eSTRs identified in this study are well
348 represented with triplet repeats that could potentially function in the same manner. The

349 substantial percentage of coding region eSTRs identified in this study may also indicate
350 previously unreported mechanisms by which they may function as cis-regulatory elements.
351 Microsatellites in 3'UTRs are assumed to influence transcript stability via AU rich repeats that
352 influence mRNA decay (Mignone et al. 2002). Mononucleotides in 3'UTRs have been proposed
353 to play a role in regulating gene expression in number of cancer-related genes (Paun et al. 2009;
354 Ziqiang et al. 2009). Rattenbacher et al.(2010) reported that GU rich elements in 3'UTRs can
355 cause mRNA decay. Collectively, these proposed mechanisms may explain how specific eSTRs
356 in 3'UTRs identified in this study may influence gene expression levels.

357 Tract length comparisons among eSTRs from the three regions did not indicate any
358 significant difference suggesting that eSTRs may expand or contract under similar selective
359 pressures irrespective of the region they are located. This finding runs contrary to previous
360 studies that report significant tract length differences in general among microsatellites located in
361 different regions. Pramod *et al.* (2014), in examining the sunflower EST database, found that
362 microsatellites in coding regions were significantly shorter than those found in UTRs suggesting
363 greater lability in UTR microsatellites.

364 We predicted that eSTRs are more likely to be found within genes that are constantly in
365 need of “tuning” to respond to changes in the environment. Previous studies on bacterial species
366 reported evidence in line with this prediction. Microsatellites were linked to the activity of some
367 hypervariable regions regulating virulence at the interface of host pathogen interactions. These
368 loci were fittingly named “contingency loci” (Moxon et al. 1994; Moxon *et al.*, 2006). To test
369 whether eSTRs are found in genes that are likely to show higher levels of evolutionary lability
370 more so than others, we searched for specific GO terms enriched within eSTR-containing
371 transcripts. The enrichment of GO terms linked to “regulation of transcription” (GO:0006355),

372 “transcription factor activity” (GO:0003700), “DNA-binding” (GO:0003677), and “positive
373 regulation of transferase activity” (GO:0051347) provide strong support for both cis- and trans-
374 regulatory roles for eSTRs. Enriched GO terms such as, “spliceosomal complex” (GO:0005681),
375 “mRNA splicing, via spliceosome” (GO:0000398), and “spliceosomal tri-snRNP complex”
376 (GO:0097526) hint at specific mechanisms in which the involvement of eSTRs may be crucial
377 (Table 1; Figure 5). Other more general GO terms identified as enriched within the eSTR-
378 containing transcripts (Supplemental Table 16) also indicate specific genes where environment
379 tracking and “tuning” may be desired.

380 Our study identified 479 transcribed microsatellites that can potentially serve as tuning
381 knobs in common sunflower. Given that our study was limited to populations across a narrow
382 latitudinal range, the number of microsatellites that show a significant effect on gene expression
383 is noteworthy. Based on these findings, we envision that the number of microsatellites that could
384 potentially alter phenotypes may be more than we could discover given the limited number of
385 microsatellites investigated in this study. Our results based on transcribed regions suggest the
386 existence of a substantial number of functional microsatellites that are potentially adaptive and
387 may indicate the need to exercise caution when using transcribed microsatellites as neutral
388 molecular markers in population genetic studies. Further, the results presented in this study
389 consistent with the proposed functionality of microsatellites indicate the potential to predict a
390 range of phenotypes based on specific microsatellite genotypes. This study provides strong
391 evidence to suggest that microsatellites can rapidly generate heritable genetic variation which
392 improves our understanding of mechanisms that can influence rapid adaptation via mutation at
393 rates far greater than typically assumed.

394

395 **METHODS**

396 *Sample collection and common garden experiment*

397 As previously described, seeds were collected from six natural populations of common
398 sunflowers from two latitudinal locations in Kansas and Oklahoma (Supplemental Table S1)
399 (Ranathunge *et al.* 2018). Scarified seed were germinated on moist filter paper in petri dishes.
400 Seedlings were transferred into 2.54 cm “cone-tainers” (Stuwe & Sons, Inc., Tangent, OR, USA)
401 arranged in a randomized design and were kept in a greenhouse under controlled conditions. At
402 the age of four weeks, young leaf tissue samples from 96 individuals representing the 6
403 populations (16 individuals from each) were collected for RNA isolation. Multiple populations
404 from each latitude were used so that the relative effect of the microsatellite allele length on gene
405 expression while controlling for other components of phenotypic variation such as variation in
406 the environment and the local genetic background.

407 *RNA-Seq and de novo transcriptome assembly*

408 RNA was isolated from 20 mg of fresh leaf tissue with Maxwell 16 LEV *simplyRNA*
409 Tissue kits (Promega, WI, USA). Isolated RNA samples were sent to HudsonAlpha Institute for
410 Biotechnology (<http://hudsonalpha.org>) for high throughput sequencing. The procedure for
411 cDNA library preparation and RNA sequencing was described in detail in Ranathunge *et al*
412 (2018). RNA sequencing (2 x 100 paired end) was carried out with Illumina HiSeq 2500
413 platform. High quality reads were obtained for 95 out of 96 individuals. The average total
414 number of reads per individual was 42.8 million. The total number of reads per individual ranged
415 from 27.6 – 88.1 million (Supplemental Table S2). RNA-Seq produced 100 bp paired end reads.
416 The quality of the reads was assessed with Fast QC software
417 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) using default settings. A deviation

418 from expected base frequencies was observed at the beginning of sequenced fragments,
419 indicating possible adapter contamination; these positions were removed from the reads by
420 trimming 14 bp from each end. Trimmed reads complying with the default quality standards of
421 the FastQC software were used in the downstream processes. The sample that produced the
422 highest total number of reads (88.1 million) which indicated the best coverage and also complied
423 with base quality standards was used to construct a reference transcriptome. The reference
424 transcriptome was built with the software modules (Inchworm, Chrysalis, and Butterfly) of the
425 Trinity program (Grabherr et al. 2011). We used the paired end option in Trinity with a minimum
426 contig size set to 200 bp and the minimum k-mer coverage set to 1. The reference transcriptome
427 consisted of 58,431 contigs. Fastq files with raw reads obtained from the remaining 94
428 individuals were aligned to the reference transcriptome with Bowtie2 (Langmead and Salzberg
429 2012). Bowtie2 produced contigs were managed using SAMtools (Li et al. 2009). SAMtools
430 produced BAM format files which were then sorted and indexed to quantify gene expression.
431 Gene expression was quantified from reads aligned to the reference transcriptome. The reads for
432 each contig were normalized by the library size of each individual, and gene expression was
433 estimated as read counts per 100 million reads.

434 ***Functional annotation***

435 A standalone BLAST (Basic Local Alignment Search Tool) search (Altschul et al. 1997)
436 was performed on the reference transcriptome against the *Helianthus annuus* protein sequence
437 database to annotate the reference transcriptome. First, protein sequence data for sunflower
438 unigenes were downloaded from the sunflower genome database at
439 <https://www.sunflowergenome.org/> and a database was created. Sequences from the reference
440 transcriptome were used as the query in the BLASTX search. An E-value cutoff of 0.0001, gap

441 open penalty score of 11, gap extension penalty score of one and minimum word size of three
442 were used as alignment parameters. BLOSUM62 was used as the matrix of choice. A best hit
443 overhang of 0.25 and maximum target sequence value of one were used to minimize the number
444 of hits for each query sequence. The results of the BLASTX search were output in tabular format
445 and the hits were further filtered based on bit score and e-value to retain the single hit with the
446 highest bit score and lowest E-value for each query sequence.

447 *Microsatellite genotyping*

448 We searched the reference transcriptome for microsatellites with Tandem Repeat Finder
449 (Benson 1999). The match, mismatch, indel and alignment scores were set at 2, 7, 7, and 50
450 respectively. The parsed output file from Tandem Repeat Finder and BAM format files from the
451 95 individuals were used as input for RepeatSeq (Highnam et al. 2013) for genotyping
452 microsatellites. RepeatSeq uses a Bayesian approach for making variant calls on reads generated
453 with short-read sequencing technology. RepeatSeq's default settings were used to genotype
454 transcribed microsatellites across all individuals. The output from RepeatSeq was initially used
455 to assess the motif size and type frequencies within the list of genotyped microsatellites. Motif
456 types were standardized with custom PERL scripts as per the motif matrix mentioned in Kofler
457 et al (2007). Output files from RepeatSeq were also used to extract position of the microsatellite
458 in the transcript.

459 Reads that failed to generate genotypes were removed from the downstream processes.
460 Raw expression estimates based on read counts were obtained for each allele for each individual
461 with data corresponding to the major component, allele genotype, reference genotype, and motif.
462 We noticed that RepeatSeq tends to erroneously call certain genotypes as heterozygous by
463 identifying underrepresented error reads as alleles. To solve this issue, we identified a genotype

464 as heterozygous only when the second most represented allele was called with at least 25% of the
465 frequency of the dominant allele. In the majority of cases with more than two alleles present,
466 many were at very low frequency and were removed by this method.

467 ***Effect of microsatellite allele length on gene expression***

468 We examined the effect of microsatellite allele length on gene expression (log-
469 transformed) by employing an analysis of covariance (ANCOVA) while controlling for
470 population and allele-by-population interaction effects. Linear, quadratic, and cubic regression
471 were all attempted because previous *in-vivo* experiments have identified non-linear relationships
472 between microsatellite allele length and gene expression (Vinces et al. 2009). A p-value of 0.05,
473 and q-value cutoff of 0.05 that corrects for multiple comparisons (Storey 2002), were used to
474 identify microsatellites where allele length correlates with gene expression (hereinafter referred
475 to as eSTRs - a term borrowed from Gymrek et al. 2015 which refers to significant expression
476 short tandem repeats). All statistical analyses were performed in R version 3.3.3 (R Core Team
477 2017). To identify eSTRs, we first employed a conservative approach by conducting ANCOVAs
478 with unfiltered microsatellite allele lengths and allele specific gene expression estimates. This
479 first analysis included irregular allele sizes inconsistent with the length of the repeat motif. We
480 then removed these alleles from the analyses and reran the ANCOVA to estimate the
481 microsatellite allele length effect on allele-specific gene expression.

482 Further, information from the RepeatSeq output was used to calculate motif size and type
483 frequencies within eSTRs to assess whether microsatellites of specific motif size and type classes
484 are more likely to function as tuning knobs. We used the BLASTX output for the reference
485 transcriptome to extract best hits for the eSTR-containing transcripts. Information on the reading
486 frame, “query start” and “query end” positions from the BLASTX output along with eSTR

487 position data from the RepeatSeq output were used to identify whether an eSTR was located
488 within the 5' UTR (Untranslated region), coding region, or the 3'UTR. To test whether
489 mechanisms underlying expansion and contraction of eSTRs are likely to be different based on
490 the location of eSTRs and/or their respective motif sizes, we conducted Kruskal-Wallis (KW)
491 tests on eSTR tract lengths. The KW tests statistically compared eSTR tract lengths among the
492 three regions, 5'UTR, coding region and 3'UTR, and among eSTRs of different motif sizes.

493 ***Validation of gene expression estimates by Real Time PCR (qPCR)***

494 We conducted qPCR analysis on seven eSTRs selected based on the magnitude of allele
495 length effect on gene expression to validate RNA-Seq based gene expression estimates. The
496 isolated RNA samples from 48 of the individuals used in the RNA-Seq experiment were
497 converted to cDNA with High Capacity cDNA Reverse Transcription kit with RNase inhibitor
498 (Applied Biosystems, Foster City, California, USA). Sequence data were obtained for seven
499 eSTR-containing transcripts from the *de novo* transcriptome. TaqMan gene expression assays
500 were designed for the seven selected loci using Primer Express version 3.0 (Applied Biosystems,
501 Foster City, California, USA) (Supplemental Table S3). TaqMan assays for two constitutively
502 expressed genes in sunflower, actin and ubiquitin, previously designed by Pramod et al. (2012)
503 were used as controls. Protocols for generating standard curves and validating RNA-Seq derived
504 gene expression estimates are given in Supplemental Note 1 and Supplemental Table S4.

505 ***Validation of RNA-Seq derived microsatellite genotypes with fragment analysis***

506 Microsatellite genotypes obtained from RepeatSeq based on RNA-Seq data were
507 validated with traditional fragment analysis on eSTRs across the 95 individuals used in the RNA-
508 Seq experiment. DNA was isolated with Maxwell 16 tissue DNA purification kit (Promega, WI,
509 USA) from dried leaf tissue from the 95 plants and primers were designed for three eSTR loci

510 with Primer 3 development software (Untergasser et al. 2012). Forward primers were
511 fluorescently labeled (Supplemental Table S5). PCR was performed with touchdown PCR
512 conditions (Don et al. 1991) (Supplemental Note 2). Fragment analysis was performed on ABI
513 3730 capillary sequencers (Applied Biosystems, USA) at the Arizona State University DNA
514 laboratory with LIZ-500 as the size standard (GeneScan – 500 LIZ Size Standard – Applied
515 Biosystems, USA). Microsatellite genotypes were visually scored using Peak Scanner version
516 1.0 (Applied Biosystems, USA).

517 RepeatSeq's genotype calling method only targets the microsatellite repeat regions and,
518 at times, the immediate flanking regions whereas fragment analysis derived microsatellite allele
519 lengths include more of the flanking sequences that are largely non-repetitive. We altered the
520 allele lengths obtained from fragment analysis by subtracting the length of the non-repetitive
521 portion of the amplicon and used this key to predict the RepeatSeq genotypes for each
522 microsatellite locus to assess the concordance between the two types of data.

523 ***Gene Ontology (GO) enrichment analysis***

524 Protein sequences corresponding to the best hits for the reference transcriptome obtained
525 from the BLASTX search were retrieved. The sequences were used to conduct a Gene Ontology
526 (GO) analysis with Blast2GO (Conesa et al. 2005) to identify GO terms associated with all genes
527 in the reference transcriptome. A BLAST search was conducted against the *Arabidopsis thaliana*
528 protein sequence database with the blastp option. An E-value cutoff of 0.001, a minimum
529 number of BLAST hits at 20, a HSP length cutoff of 33, and a word size of 3 were used as
530 settings in the blastp search. The BLAST hits were then mapped and annotated with default
531 settings to identify GO terms under three categories, namely, biological processes, molecular
532 function, and cellular component. To identify specific GO terms enriched within eSTR-

533 containing transcripts, Fisher's exact test was performed with GO terms associated with all *H.*
534 *annuus* genes in the reference transcriptome as the background list and *H.annuus* gene IDs
535 corresponding to eSTR-containing transcripts as the test set. A False Discovery Rate (FDR) of
536 0.05 was used as the significance threshold to identify overrepresented GO terms within the
537 eSTR-containing transcripts. General GO terms (parent terms) were removed to retrieve the most
538 specific terms using "reduce to most specific" option in Blast2GO. The REVIGO (Supek et al.
539 2011) tool was also used to reduce the enriched GO terms to the most specific terms for
540 visualization based on the uniqueness and dispensability scores calculated by the program.

541 **DATA ACCESS**

542 All sequence data have been deposited at the National Center for Biotechnology Information
543 short read archive under project PRJNA408292.

544 **ACKNOWLEDGEMENTS**

545 This work was supported by the National Science Foundation grants, MCB- 1158521 to
546 M.E.W and EPS-0903787 to A.D.P, and the Department of Biological Sciences at Mississippi
547 State University. The authors wish to thank Lisa E. Wallace for specimen voucher preparation,
548 Jessica L. Martin Judson for help with sample collection, and Brian S. Baldwin and Jesse I.
549 Morrison for assistance with the common garden experiment.

550

551

552

553

554

555 **REFERENCES**

- 556 Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. 1997. Gapped
557 BLAST and PSI-BLAST: A new generation of protein database search programs. *Nucleic*
558 *Acids Res* **25**: 3389–3402.
- 559 Andrew SE, Goldberg YP, Kremer B, Telenius H, Theilmann J, Adam S, Starr E, Squitieri F, Lin
560 B, Kalchman MA, et al. 1993. The relationship between trinucleotide (CAG) repeat length
561 and clinical features of Huntington’s disease. *Nat Genet* **4**: 398–403.
- 562 Bachtrog D, Weiss S, Zangerl B, Brem G, Schlötterer C. 1999. Distribution of dinucleotide
563 microsatellites in the *Drosophila melanogaster* genome. *Mol Biol Evol* **16**: 602–610.
- 564 Barton NH. 1990. Pleiotropic models of quantitative variation. *Genetics* **124**: 773–782.
- 565 Benson G. 1999. Tandem repeats finder: A program to analyze DNA sequences. *Nucleic Acids*
566 *Res* **27**: 573–580.
- 567 Blackman BK, Michaels SD, Rieseberg LH. 2011. Connecting the sun to flowering in sunflower
568 adaptation. *Mol Ecol* **20**: 3503–3512.
- 569 Chung M-Y, Ranum LPW, Duvick LA, Servadio A, Zoghbi HY, Orr HT. 1993. Evidence for a
570 mechanism predisposing to intergenerational CAG repeat instability in spinocerebellar
571 ataxia type I. *Nat Genet* **5**: 254–258.
- 572 Conesa A, Götz S, García-Gómez JM, Terol J, Talón M, Robles M. 2005. Blast2GO: A universal
573 tool for annotation, visualization and analysis in functional genomics research.
574 *Bioinformatics* **21**: 3674–3676.
- 575 Contente A, Dittmer A, Koch MC, Roth J, Dobbelstein M. 2002. A polymorphic microsatellite
576 that mediates induction of PIG3 by p53. *Nat Genet* **30**: 315–320.
- 577 Costa R, Peixoto AA, Thackeray JR, Dalgleish R, Kyriacou CP. 1991. Length polymorphism in
578 the threonine-glycine-encoding repeat region of the period gene in *Drosophila*. *J Mol Evol*
579 **32**: 238–246.
- 580 Don RH, Cox PT, Wainwright BJ, Baker K, Mattick JS. 1991. “Touchdown” PCR to circumvent
581 spurious priming during gene amplification. *Nucleic Acids Res* **19**: 4008.
- 582 Dudley JW, Lambert RJ. 1992. Ninety generations of selection for oil and protein in maize.
583 *Maydica*.
- 584 Eichler EE, Holden JJA, Popovich BW, Reiss AL, Snow K, Thibodeau SN, Richards CS, Ward
585 PA, Nelson DL. 1994. Length of uninterrupted CGG repeats determines instability in the
586 FMR1 gene. *Nat Genet* **8**: 88–94.
- 587 Ellegren H. 2004. Microsatellites: Simple sequences with complex evolution. *Nat Rev Genet* **5**:
588 435–445.
- 589 Fahima T, Röder MS, Wendehake K, Kirzhner VM, Nevo E. 2002. Microsatellite polymorphism
590 in natural populations of wild emmer wheat, *Triticum dicoccoides*, in Israel. *Theor Appl*
591 *Genet* **104**: 17–29.

- 592 Fondon JW, Garner HR. 2004. Molecular origins of rapid and continuous morphological
593 evolution. *Proc Natl Acad Sci U S A* **101**: 18058–18063.
- 594 Gemayel R, Chavali S, Pougach K, Legendre M, Zhu B, Boeynaems S, van der Zande E,
595 Gevaert K, Rousseau F, Schymkowitz J, et al. 2015. Variable glutamine-rich repeats
596 modulate transcription factor activity. *Mol Cell* **59**: 615–627.
- 597 Gemayel R, Vinces MD, Legendre M, Verstrepen KJ. 2010. Variable tandem repeats accelerate
598 evolution of coding and regulatory sequences. *Annu Rev Genet* **44**: 445–477.
- 599 Golubov A, Yao Y, Maheshwari P, Bilichak A, Boyko A, Belzile F, Kovalchuk I. 2010.
600 Microsatellite instability in *Arabidopsis* increases with plant development. *Plant Physiol*
601 **154**: 1415–1427.
- 602 Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,
603 Raychowdhury R, Zeng Q, et al. 2011. Full-length transcriptome assembly from RNA-Seq
604 data without a reference genome. *Nat Biotechnol* **29**: 644–652.
- 605 Gymrek M, Willems T, Guilmatre A, Zeng H, Markus B, Georgiev S, Daly MJ, Price AL,
606 Pritchard JK, Sharp AJ, et al. 2015. Abundant contribution of short tandem repeats to gene
607 expression variation in humans. *Nat Genet* **48**: 22–29.
- 608 Hammock EAD, Young LJ. 2005. Genetics: Microsatellite instability generates diversity in brain
609 and sociobehavioral traits. *Science (80-)* **308**: 1630–1634.
- 610 Heiser CB, Smith DM, Clevenger SB, Martin WC. 1969. The north american sunflowers
611 (*Helianthus*). *Mem Torrey Bot Club* **22**: 1–218.
- 612 Highnam G, Franck C, Martin A, Stephens C, Puthige A, Mittelman D. 2013. Accurate human
613 microsatellite genotypes from high-throughput resequencing data using informed error
614 profiles. *Nucleic Acids Res* **41**.
- 615 Hopkins WD, Donaldson ZR, Young LJ. 2012. A polymorphic indel containing the RS3
616 microsatellite in the 5' flanking region of the vasopressin V1a receptor gene is associated
617 with chimpanzee (*Pan troglodytes*) personality. *Genes, Brain Behav* **11**: 552–558.
- 618 Jarne P, Lagoda PJJ. 1996. Microsatellites, from molecules to populations and back. *Trends*
619 *Ecol Evol* **11**: 424–429.
- 620 Kashi Y, King DG. 2006. Simple sequence repeats as advantageous mutators in evolution.
621 *Trends Genet* **22**: 253–259.
- 622 King DG, Soller M, Kashi Y. 1997. Evolutionary tuning knobs. *Endeavour* **21**: 36–40.
- 623 Kofler R, Schlötterer C, Lelley T. 2007. SciRoKo: A new tool for whole genome microsatellite
624 search and investigation. *Bioinformatics* **23**: 1683–1685.
- 625 Kumar S, Bhatia S. 2016. A polymorphic (GA/CT)_n-SSR influences promoter activity of
626 Tryptophan decarboxylase gene in *Catharanthus roseus* L. Don. *Sci Rep* **6**.
- 627 Kunkel TA, Soni A. 1988. Mutagenesis by transient misalignment. *J Biol Chem* **263**: 14784–
628 14789.

- 629 Kunst CB, Leeflang EP, Iber JC, Arnheim N, Warren ST. 1997. The effect of FMR1 CGG repeat
630 interruptions on mutation frequency as measured by sperm typing. *J Med Genet* **34**: 627–
631 631.
- 632 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**:
633 357–359.
- 634 Lawson MJ, Zhang L. 2006. Distinct patterns of SSR distribution in the *Arabidopsis thaliana*
635 and rice genomes. *Genome Biol* **7**.
- 636 Lee K, Dunlap JC, Loros JJ. 2003. Roles for white collar-1 in circadian and general
637 photoperception in *Neurospora crassa*. *Genetics* **163**: 103–114.
- 638 Lee TH, Maheshri N. 2012. A regulatory role for repeated decoy transcription factor binding
639 sites in target gene expression. *Mol Syst Biol* **8**.
- 640 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R.
641 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**: 2078–2079.
- 642 Li WH. 1997. Molecular evolution Sinauer. *Sunderland, Mass.*
- 643 Li YC, Korol AB, Fahima T, Beiles A, Nevo E. 2002. Microsatellites: Genomic distribution,
644 putative functions and mutational mechanisms: A review. *Mol Ecol* **11**: 2453–2465. inter
- 645 Li YC, Korol AB, Fahima T, Nevo E. 2004. Microsatellites within genes: Structure, function,
646 and evolution. *Mol Biol Evol* **21**: 991–1007.
- 647 Linder CR. 2000. Adaptive evolution of seed oils in plants: Accounting for the biogeographic
648 distribution of saturated and unsaturated fatty acids in seed oils. *Am Nat* **156**: 442–458.
- 649 Mackay TFC. 1995. The genetic basis of quantitative variation: numbers of sensory bristles of
650 *Drosophila melanogaster* as a model system. *Trends Genet* **11**: 464–470.
- 651 Mignone F, Gissi C, Liuni S, Pesole G. 2002. Untranslated regions of mRNAs. *Genome Biol* **3**:
652 0004.1-0004.10.
- 653 Morgante M, Hanafey M, Powell W. 2002. Microsatellites are preferentially associated with
654 nonrepetitive DNA in plant genomes. *Nat Genet* **30**: 194–200.
- 655 Moxon ER, Rainey PB, Nowak MA, Lenski RE. 1994. Adaptive evolution of highly mutable
656 loci in pathogenic bacteria. *Curr Biol* **4**: 24–33.
- 657 Moxon R, Bayliss C, Hood D. 2006. Bacterial contingency loci: The role of simple sequence
658 DNA repeats in bacterial adaptation. *Annu Rev Genet* **40**: 307–333.
- 659 Nevo E, Beharav A, Meyer RC, Hackett CA, Forster BP, Russell JR, Powell W. 2005. Genomic
660 microsatellite adaptive divergence of wild barley by microclimatic stress in “Evolution
661 Canyon”, Israel. *Biol J Linn Soc* **84**: 205–224.
- 662 Orr HA. 2010. The population genetics of beneficial mutations. *Philos Trans R Soc B Biol Sci*
663 **365**: 1195–1201.
- 664

- 665 Paun BC, Cheng Y, Leggett BA, Young J, Meltzer SJ, Mori Y. 2009. Screening for
666 microsatellite instability identifies frequent 3'-Untranslated region mutation of the RB1-
667 inducible coiled-coil 1 gene in colon tumors. *PLoS One* **4**.
- 668 Pramod S, Downs KE, Welch ME. 2012. Gene expression assays for actin, ubiquitin, and three
669 microsatellite-encoding genes in *Helianthus annuus* (Asteraceae). *Am J Bot* **99**: 350–352.
- 670 Pramod S, Perkins A, Welch M. 2014. Patterns of microsatellite evolution inferred from the
671 *Helianthus annuus* (Asteraceae) transcriptome. *J Genet* **93**: 431–442.
- 672 Qin Z, Wang Y, Wang Q, Li A, Hou F, Zhang L. 2015. Evolution analysis of simple sequence
673 repeats in plant genome. *PLoS One* **10**.
- 674 Ranathunge C, Wheeler GL, Chimahusky ME, Kennedy MM, Morrison JI, Baldwin BS, Perkins
675 AD, Welch ME. 2018. Transcriptome profiles of sunflower reveal the potential role of
676 microsatellites in gene expression divergence. *Mol Ecol* **27**: 1188–1199.
- 677 R Core Team (2017) R: A language and environment for statistical computing. R Foundation for
678 Statistical Computing, Vienna, Austria URL <https://www.R-project.org>.
- 679 Rattenbacher B, Beisang D, Wiesner DL, Jeschke JC, Von Hohenberg M, St. Louis-Vlasova IA,
680 Bohjanen PR. 2010. Analysis of CUGBP1 targets identifies GU-repeat sequences that
681 mediate rapid mRNA decay. *Mol Cell Biol* **30**: 3970–3980.
- 682 Ryan PR, Raman H, Gupta S, Sasaki T, Yamamoto Y, Delhaize E. 2010. The multiple origins of
683 aluminium resistance in hexaploid wheat include *Aegilops tauschii* and more recent cis
684 mutations to TaALMT1. *Plant J* **64**: 446–455.
- 685 Sandman K, Reeve JN. 1999. Archaeal nucleosome positioning by CTG repeats. *J Bacteriol* **181**:
686 1035–1038.
- 687 Shimajiri S, Arima N, Tanimoto A, Murata Y, Hamada T, Wang KY, Sasaguri Y. 1999.
688 Shortened microsatellite d(CA)₂₁ sequence down-regulates promoter activity of matrix
689 metalloproteinase 9 gene. *FEBS Lett* **455**: 70–74.
- 690 Storey JD. 2002. A direct approach to false discovery rates. *J R Stat Soc Ser B Stat Methodol* **64**:
691 479–498.
- 692 Streelman JT, Kocher TD. 2002. Microsatellite variation associated with prolactin expression
693 and growth of salt-challenged tilapia. *Physiol Genomics* **2002**: 1–4.
- 694 Supek F, Bošnjak M, Škunca N, Šmuc T. 2011. Revigo summarizes and visualizes long lists of
695 gene ontology terms. *PLoS One* **6**.
- 696 Tautz D, Renz M. 1984. Simple sequences are ubiquitous repetitive components of eukaryotic
697 genomes. *Nucleic Acids Res* **12**: 4127–4138.
- 698 Tóth G, Gáspári Z, Jurka J. 2000. Microsatellites in different eukaryotic genomes: Surveys and
699 analysis. *Genome Res* **10**: 967–981.
- 700 Toutenhoofd SL, Garcia F, Zacharias DA, Wilson RA, Strehler EE. 1998. Minimum CAG repeat
701 in the human calmodulin-1 gene 5' untranslated region is required for full expression.
702 *Biochim Biophys Acta - Gene Struct Expr* **1398**: 315–320.

- 703 Trifonov EN. 2004. Tuning function of tandemly repeating sequences: a molecular device for
704 fast adaptation. In *Evolutionary theory and processes: Modern horizons*, pp. 115–138,
705 Springer.
- 706 Untergasser A, Cutcutache I, Koressaar T, Ye J, Faircloth BC, Remm M, Rozen SG. 2012.
707 Primer3-new capabilities and interfaces. *Nucleic Acids Res* **40**.
- 708 Verkerk AJMH, Pieretti M, Sutcliffe JS, Fu YH, Kuhl DPA, Pizzuti A, Reiner O, Richards S,
709 Victoria MF, Zhang F, et al. 1991. Identification of a gene (FMR-1) containing a CGG
710 repeat coincident with a breakpoint cluster region exhibiting length variation in fragile X
711 syndrome. *Cell* **65**: 905–914.
- 712 Verstrepen KJ, Jansen A, Lewitter F, Fink GR. 2005. Intragenic tandem repeats generate
713 functional variability. *Nat Genet* **37**: 986–990.
- 714 Vences MD, Legendre M, Caldara M, Hagihara M, Verstrepen KJ. 2009. Unstable tandem
715 repeats in promoters confer transcriptional evolvability. *Science (80-)* **324**: 1213–1216.
- 716 Vogt P. 1990. Potential genetic functions of tandem repeated DNA sequence blocks in the
717 human genome are based on a highly conserved “chromatin folding code.” *Hum Genet* **84**:
718 301–336.
- 719 Wang YH. 2007. Chromatin structure of repeating CTG/CAG and CGG/CCG sequences in
720 human disease. *Front Biosci* **12**: 4731–4741.
- 721 Yoo BH, Nicholas FW, Rathie KA. 1980. Long-term selection for a quantitative character in
722 large replicate populations of *Drosophila melanogaster* - Part 4: Relaxed and reverse
723 selection. *Theor Appl Genet* **57**: 113–117.
- 724 Ziqiang Y, Joongho S, Wilson A, Goel S, Ling YH, Ahmed N, Dopeso H, Jhaver M, Nasser S,
725 Montagna C, et al. 2009. An A13 repeat within the 3'-untranslated region of epidermal
726 growth factor receptor (EGFR) is frequently mutated in microsatellite instability colon
727 cancers and is associated with increased EGFR expression. *Cancer Res* **69**: 7811–7818.
- 728