1 **Characterization of duplicate gene evolution in the recent**

2 **natural allopolyploid *Tragopogon miscellus* by next-generation**

3 **sequencing and Sequenom iPLEX MassARRAY genotyping**

4

5 Richard J. A. Buggs[1,2,][*], Srikar Chamala[1, 3,][*], Wei Wu[4], Lu Gao[4], Gregory D. May[5],

6 Patrick S. Schnable[4], Douglas E. Soltis[1, 3], Pamela S. Soltis[2,3], W. Brad Barbazuk[1, 3]

7

8 [1]Department of Biology, University of Florida, Gainesville, Florida 32611.

9 [2]Florida Museum of Natural History, University of Florida, Gainesville, Florida 32611.

10 [3]Genetics Institute, University of Florida, Gainesville, Florida 32610.

11 [4]Center for Plant Genomics, Iowa State University, Ames, Iowa 50011.

12 [5]National Center for Genome Resources, Santa Fe, New Mexico 87505.

13 *These authors contributed equally to this manuscript.

14

15 Corresponding author: Dr W. Brad Barbazuk, Department of Biology, University of

16 Florida, Gainesville, Florida 32611, USA.

17 Email: bbarbazuk@ufl.edu

18 Telephone: 352-273-8624

19 Key words: polyploidy, whole genome duplication, *Tragopogon miscellus*, homoeolog

20 evolution, pyrosequencing

21 Running title: Homoeolog loss in a natural allopolyploid

22 Word count: 5813

23

# Abstract

*Tragopogon miscellus* (Asteraceae) is an evolutionary model for the study of natural allopolyploidy, but until now has been under-resourced as a genetic model. Using 454 and Illumina expressed sequence tag sequencing of the parental diploid species of *T. miscellus*, we identified 7782 single nucleotide polymorphisms that differ between the two progenitor genomes present in this allotetraploid. Validation of a sample of 98 of these SNPs in genomic DNA using Sequenom MassARRAY iPlex genotyping confirmed 92 SNP markers at the genomic level that were diagnostic for the two parental genomes. In a transcriptome profile of 2989 SNPs in a single *T. miscellus* leaf, using Illumina sequencing, 69% of SNPs showed approximately equal expression of both homeologs (duplicate homologous genes derived from different parents), 22% showed apparent differential expression, and 8.5% showed apparent silencing of one homeolog in *T. miscellus*. The majority of cases of homeolog silencing involved the *T. dubius* SNP homeolog (164/254; 65%) rather than the *T. pratensis* homeolog (90/254). Sequenom analysis of genomic DNA showed that in a sample of 27 of the homeologs showing apparent silencing, 23 (85%) were due to genomic homeolog loss. These methods could be applied to any organism, allowing efficient and cost-effective generation of genetic markers.

# Introduction

44

45  Many natural and domesticated plant species are hybrids which have undergone whole-

46  genome duplication. This condition, known as allopolyploidy (Kihara & Ono 1927), may

47  have large effects on both the ecology (e.g. Stebbins 1942; Buggs & Pannell 2007) and

48  evolution (Soltis & Soltis 1999; Adams & Wendel 2005) of a lineage. Genome evolution

49  of allopolyploids has been extensively studied in crop species such as cotton (Adams &

50  Wendel 2004; Udall & Wendel 2006), wheat (Feldman *et al.* 1997; Levy & Feldman

51  2004; Dong *et al.* 2005; Bottley *et al.* 2006), soybean (Joly *et al.* 2004), and tobacco (Lim

52  *et al.* 2004; Petit *et al.* 2007), as well as genetic models such as *Arabidopsis* (Chen *et al.*

53  2004; Chen *et al.* 2008). These studies demonstrate dynamic patterns of evolution, but

54  have limitations due to uncertainties about the precise history and ecological context of

55  the lineages. Furthermore, they cannot provide insights into the early stages of polyploid

56  evolution in nature.  It is therefore difficult to know whether certain evolutionary changes

57  took place in the progenitor diploids, upon allopolyploidization, or in the subsequent

58  generations.

59

60  A need therefore exists for natural allopolyploid model organisms with a known history

61  and ecological context (Soltis *et al.* 2004b; Buggs 2008). A handful of species have been

62  identified for this purpose, such as *Senecio cambrensis* (Hegarty *et al.* 2005), *Spartina*

63  *anglica* (Ainouche *et al.* 2004), *Tragopogon mirus* and *T. miscellus* (Soltis *et al.* 2004a).

64  *Tragopogon miscellus* is a particularly tractable evolutionary model for the study of the

65 early generations of allopolyploidy. Its origin can be accurately dated to about 80 years

66 ago (Ownbey 1950; Soltis *et al.* 2004a). The parental diploid species are known and still

67 coexist with their allopolyploid derivative; both reciprocal crosses of the parents exist in

68 natural populations, and at least one of them appears to have originated multiple times

69 (Novak *et al.* 1991; Soltis *et al.* 1995; Symonds *et al.* 2009). *Tragopogon miscellus* is a

70 textbook example of allopolyploid speciation (e.g. Judd *et al.* 2007; Sadava *et al.* 2008).

71

72 Unlike the crop species that have been used to study allopolyploid evolution, the natural

73 allopolyploid evolutionary model systems are under-resourced as genetic models. To

74 date, the best resourced is *S. cambrensis* for which cDNA microarrays have been made to

75 study gene expression (Hegarty *et al.* 2005; Hegarty *et al.* 2006). Until now resources for

76 *T. miscellus* have consisted of DNA sequence tags for only 23 duplicate gene pairs (Tate

77 *et al.* 2006; Buggs *et al.* 2009; Tate *et al.* 2009a), a handful of phylogenetic markers

78 (Mavrodiev *et al.* 2005), and 2000 uncharacterized Sanger ESTs (J. Koh, J. Tate, D.

79 Soltis and P. Soltis, unpubl. data). This paucity of sequence data contrasts with the

80 usefulness of *T. miscellus* as an evolutionary model.

81

82 One key issue in the evolution of allopolyploids is the fate of duplicated genes. Duplicate

83 gene evolution is important for understanding the evolution of the allopolyploids

84 themselves, and may allow for more general statements about the evolution of duplicated

85 genes in non-polyploid organisms. Natural allopolyploid models present systems

86 containing a whole genome's worth of duplicated genes of identical and known age.

87 Duplicated genes may have a variety of evolutionary fates: non-functionalization, sub-

88    functionalization and neo-functionalization (Lynch & Conery 2000). Several studies have

89    examined the evolution of homeologs (genes duplicated by whole-genome duplication) in

90    allopolyploids. Studies in crop species have shown homeolog loss (e.g. Song *et al.* 1995;

91    Kashkush *et al.* 2002) and patterns of homeolog expression suggestive of

92    subfunctionalization (e.g. Adams *et al.* 2003; Flagel *et al.* 2008).

93

94    In natural models, our knowledge of homeolog evolution is limited. In the *S. cambrensis*

95    cDNA microarray, the oligo-nucleotides used did not distinguish between homeologs:

96    measures of gene expression were the total expression of both homeologs. In *T.*

97    *miscellus*, loss and silencing of homeologs occurred in the early generations of

98    allopolyploidy (Tate *et al.* 2006; Buggs *et al.* 2009; Tate *et al.* 2009a) based on analysis

99    of only 20 homeolog pairs using PCR-based methods. New surveys are needed that will

100   move us from a gene-by-gene approach to a genomic level. This requires a dramatic

101   increase in the genomic resources available for plants that are good evolutionary models

102   but not genetic models. We wished to develop a protocol that would produce a large

103   number of homeolog-specific markers in *T. miscellus* at minimal time and expense,

104   allowing us to assess homeologous gene loss and silencing.

105

106   Sequencing of cDNA or expressed sequence tags (EST) provides a rapid method for gene

107   discovery and can be used to identify transcripts associated with specific biological

108   processes. As such, it is often a first step in the genomic characterization of an organism.

109   Variation in ESTs can be characterized by single nucleotide polymorphisms (SNPs),

110   which are single-base differences between haplotypes. Transcript-associated SNPs can be

111   used to develop allele-specific assays for the examination of *cis*-regulatory variation

112   within a species (Guo *et al.* 2004; Stupar & Springer 2006) and may provide a rapid

113   means to investigate differential expression and gene gain/loss within polyploids. EST

114   collections and SNP discovery rely on DNA sequencing, which until recently was

115   prohibitively costly for most evolutionary studies.

116

117   Recent advances in high-throughput sequencing technology provide rapid and cost-

118   effective means to generate sequence data (Stupar & Springer 2006; Ellegren 2008;

119   Hudson 2008).  This new paradigm, termed flow-cell sequencing (reviewed in Holt &

120   Jones 2008), consists of stepwise determination of DNA sequence by iterative cycles of

121   nucleotide extensions done in parallel on huge numbers of clonally amplified template

122   molecules. This massively parallel approach enables DNA sequence to be acquired at

123   extremely high depths of coverage in less time and for less cost than traditional

124   sequencing. The 454-FLX produces 200,000 sequences per run with ~200-300 bp lengths

125   (100 Mb). With new Titanium reagents, this can be increased to over 1 million sequences

126   with ~350-400 bp read lengths (400-600 Mb per run). In contrast, the Illumina Genome

127   Analyzer (GA) II DNA sequencing instrument can produce >80 million sequences, each

128   of which is 36bp in length (> 2 Gb). Short read lengths can confound assembly and

129   alignment programs, but the reduction in read length vs. increased depth of coverage is an

130   acceptable trade-off for many re-sequencing applications such as transcript expression

131   profiling (Eveland *et al.* 2008), *in vivo* DNA binding site detection (Johnson *et al.* 2007)

132   and polymorphism detection (Barbazuk *et al.* 2007; Novaes *et al.* 2008; Van Tassell *et al.*

133   2008). In the latter application, a high volume of short reads is very powerful in

134    discriminating sequence variants, enabling reliable SNP discovery, so long as each read is

135    long enough and accurate enough to align uniquely to the reference sequences.

136

137    To permit gene discovery and genomic tool development in species with few genomic

138    resources, we designed a hybrid sequencing approach. In this approach, the Roche 454

139    sequencer is first used to generate transcriptome or genomic sequences that can be

140    assembled and used as reference sequences (as in, e.g. Novaes *et al.* 2008). We then use

141    this reference for subsequent alignment of Illumina short reads. This method gains

142    maximum leverage from the longer read lengths of 454 sequencing and the deeper

143    coverage of Illumina. Assembling 454 sequence reads is less problematic than Illumina

144    reads, making it the high-throughput sequencing method of choice for species with few

145    genomic resources, and it is particularly useful in transcriptome characterization (Cheung

146    *et al.* 2006; Emrich *et al.* 2007; Cheung *et al.* 2008; Novaes *et al.* 2008). The 454

147    assemblies can therefore be used for gene annotation and the Illumina sequences used to

148    identify SNPs and examine relative expression differences.

149

150    Once SNPs have been identified, a highly efficient way to validate them and carry out

151    large-scale surveys of their frequencies is the Sequenom MassARRAY iPLEX genotyping

152    platform (Gabriel *et al.* 2009). In this method, a short section of DNA containing a SNP

153    is amplified from an individual by PCR. This is followed by a high-fidelity single-base

154    primer extension reaction over the SNP being assayed, using nucleotides of modified

155    mass. The different alleles therefore produce oligonucleotides with mass differences that

156    can be detected using highly accurate Matrix-Assisted Laser Desorption/Ionization Time-

157 Of-Flight (MALDI-TOF) mass spectrometry. Up to 40 different SNPs can be multiplexed

158 in one assay if primers are designed by custom software to give unique mass ranges for

159 each SNP. This method is especially suited for detecting homeologs which differ in only

160 a few SNPs as, unlike microarrays which rely on hybridization of oligonucleotides, it

161 detects differences by single-nucleotide extension over SNPs.

162

163 In this paper we demonstrate the utility of hybrid next-generation sequencing and

164 Sequenom genotyping for the study of homeolog evolution in *T. miscellus*. We report the

165 transcriptome characterization of *T. dubius*, one of the diploid progenitors of *T. miscellus*,

166 with 454 sequencing and the subsequent discovery of over 24,000 SNPs between *T.*

167 *dubius* and the other parental diploid species, *T. pratensis*, using Illumina sequencing.

168 We validated a subset of 98 SNPs that represent homeolog pairs in *T. miscellus* at the

169 genomic level using Sequenom MassARRAY iPLEX genotyping. In addition, expression

170 profiling of a *T. miscellus* individual using Illumina sequencing was performed. We

171 assessed the utility of this profile for the selection of candidate genes for the investigation

172 of loss from the genome. These methods could be applied to any organism, allowing

173 efficient and cost-effective generation of genetic markers.

174

175

# 176 **Materials and methods**

177 Seeds were collected from natural populations of allotetraploid *T. miscellus* (Soltis and

178 Soltis collection number 2671) and its diploid parent species, *T. dubius* (collection no.

179    2674) and *T. pratensis* (collection no. 2672), in Oakesdale, WA. The three species grow

180    in sympatry in this location, and this fact, together with microsatellite data (Symonds *et*

181    *al.* 2009), suggest that the diploid populations were the source of the progenitors of the

182    allotetraploid population. These seeds were germinated and grown in an air-conditioned

183    greenhouse with supplementary lighting at the University of Florida (Gainesville, FL,

184    USA). *T. miscellus* from Oakesdale is the short-liguled form, with *T. pratensis* as the

185    maternal parent (Soltis & Soltis 1989; Soltis *et al.* 1995).

186

187    RNA was extracted from leaf tissue of three individuals from Oakesdale: *T. dubius* 2674-

188    4 (ID no. 3911), *T. pratensis* 2671-1 (ID no. 3912), and *T. miscellus* 2671-1 (ID no.

189    3912). Basal leaf tissue from each plant was flash frozen and ground in liquid nitrogen

190    using a pestle and mortar. RNA extractions were performed following a portion of the

191    CTAB DNA extraction protocol (Doyle & Doyle 1987) and subsequent use of the

192    RNeasy Plant Mini Kit (Qiagen, Stanford, CA, USA) with on-column DNase digestion.

193    This method was originally developed for the successful extraction of RNA from

194    *Amborella* and *Nuphar* (Kim *et al.* 2004) and copes well with the latex produced by

195    *Tragopogon* photosynthetic tissue. This was followed by an RNA cleanup using the

196    protocol of the RNeasy Plant Mini Kit. These extractions were quality-checked using the

197    Agilent 2100 Bioanalyzer (Agilent Technologies, Palo Alto, CA, USA).

198

199    **454 EST sequencing and processing**

200

201    Using the *T. dubius* RNA, a normalized cDNA library was produced via the following

202    method. The Evrogen MINT cDNA synthesis kit (Evrogen, Moscow, Russia) was used to

203    produce double-stranded cDNA following the manufacturer's protocol. This cDNA was

204    cleaned using the Wizard® SV Gel and PCR Clean-Up System (Promega, Madison, WI,

205    USA). The Evrogen TRIMMER cDNA normalization kit (Evrogen, Moscow, Russia)

206    was used to normalize and amplify the cDNA library, following the manufacturer's

207    instructions. In the normalization step, a 0.5 dilution of the duplex-specific nuclease

208    (DSN) was found to be optimal. In the amplification step, 12 cycles were found to be

209    optimal. The resulting normalized library was used for 454 sequencing.

210

211     454 sequencing was performed as described in the supplementary material and methods

212    to Margulies *et al.* (2005) with slight modifications as specified by 454 Life Sciences.

213    Briefly, cDNA was sheared by nebulization to a size range of 300 to 800 bp. DNA

214    fragment ends were repaired and phosphorylated using T4 DNA polymerase and T4

215    polynucleotide kinase. Adaptor oligonucleotides "A" and "B" supplied with the 454 Life

216    Sciences sequencing reagent kit were ligated to the DNA fragments using T4 DNA

217    ligase. Purified DNA fragments were hybridized to DNA capture beads and clonally

218    amplified by emulsion PCR (emPCR). DNA capture beads containing amplified DNA

219    were deposited on a 70 x 75 mm PicoTiter plate, and DNA sequences determined using

220    the GS-FLX instrument. This resulted in 822,594 EST sequences.  The *T. dubius* 454

221    EST sequences were assembled with the Newbler assembler, a part of the software

222    package distributed with 454 sequencing machines. Newbler is an assembler that takes

223    into account the specifics of pyrosequencing errors to generate accurate contigs

224     (Chaisson & Pevzner 2008). Our assembly used the default directives and a vector

225     trimming database including the Evrogen primer and 454 adapter sequences.

226

227     **Comparisons of 454 ESTs to public sequence database (Annotation)**

228

229     Assembled and annotated contig EST assemblies and singletons were obtained from the

230     curated Gene Indices Project (Quackenbush *et al.* 2000);

231     http://compbio.dfci.harvard.edu/tgi/) from three other species in the Asteraceae: *Lactuca*

232     *sativa* (ver. 3.0), *Lactuca serriola* (ver. 1.0) and *Helianthus annus* (ver. 5.0).  These

233     sequences were pooled, formatted into a blastable database and aligned to the *T. dubius*

234     454 EST assemblies with WU-TBLASTX (ver 2.0), which translates both the query and

235     subject sequences in all 6 potential reading frames prior to alignment, to identify the top

236     hit for each *T. dubius* contig (pval <= 1e-05 and <= 1e-10). The *T. dubius* 454 EST

237     contigs were also BLASTX-aligned to *Arabidopsis* CDS sequences (TAIR v. 8) because

238     *Arabidopsis* represents the best curated plant genome available. Top hits for each *T.*

239     *dubius* contig to the *Arabidopsis* protein set were identified (pval <= 1e-05 and <= 1e-

240     10).  Similarity search results are summarized in Table 1.

241

242     **Illumina sequencing**

243

244     The RNA extractions from *T. dubius* 2674-4 (ID no. 3911), *T. pratensis* 2672-5 (ID no.

245     3913), and *T. miscellus* 2671-1 (ID no. 3912) were used for Illumina sequencing. Poly

246     A+ RNA was isolated from total RNA through two rounds of oligo-dT selection

247 (Dynabeads mRNA Purification Kit, Invitrogen Inc., Santa Clara, CA, USA). The mRNA

248 was annealed to high concentrations of random hexamers and reverse transcribed.

249 Following second strand synthesis, end repair, and A-tailing, adapters complementary to

250 sequencing primers were ligated to cDNA fragments (mRNA-Seq Sample Prep Kit,

251 Illumina, San Diego, CA, USA). Resultant cDNA libraries were size fractionated on

252 agarose gels, and 250 bp fragments were excised and amplified by 15 cycles of

253 polymerase chain reaction.   Resultant libraries were quality assessed using a Bioanalyzer

254 2100 and sequenced for 36 cycles on an Illumina GA II DNA sequencing instrument

255 using standard procedures.

256

257 **SNP discovery**

258

259 All Illumina reads from the *T. dubius* and *T. pratensis* parents and the *T. miscellus*

260 allotetraploid were labeled with species identifiers, pooled and aligned to the *T. dubius*

261 454 FLX contigs with the MosaikAligner package (Hillier *et al.* 2008) using the

262 following MosaikAligner parameters: -a (alignment algorithm) all; -p (CPUs used) 8;  -

263 mm (maximum mismatch) 2 in a preliminary analysis and 1 in a final analysis;  -m

264 (alignment mode) unique; -hs (hash size) 15; -mhp (maximum number of hash positions

265 to use) 100.  These alignment parameters ensured that each Illumina sequence aligned to

266 a unique position within the 454 *T. dubius* EST assembly reference sequences and with

267 no more than one base-pair mismatch in the final analysis. Illumina reads that did not

268 align with the 454 contigs under these stringent conditions were discarded from the

269 analysis.

270

271    SNPs were identified within the alignments with the GigaBayes package

272    (http://bioinformatics.bc.edu/marthlab/GigaBayes). GigaBayes is a reimplementation of

273    the PolyBayes (Marth *et al.* 1999) SNP discovery tool that has been optimized for next-

274    generation sequences. Arguments to GigaBayes were: --D (pairwise nucleotide diversity)

275    0.001; --ploidy (sample ploidy) diploid; --sample (sequence source) multiple; --anchor; --

276    algorithm banded; --CAL (minimum overall allele coverage) 3; --QRL (minimum base

277    quality value) 20.  Custom PERL scripts were written to automate the SNP discovery

278    process on all alignments to reference contigs and to parse the GigaBayes output files

279    (GFF), which contain the site identification of each SNP, its representation within each of

280    the three *Tragopogon* species (*T. dubius, T. pratensis* and *T. miscellus*) and its allele

281    usage.

282

283    Any site where both the *T. pratensis* and *T. dubius* homeologs were evidenced in the *T.*

284    *miscellus* data was flagged as a suitable SNP for the study of homeolog loss in *T.*

285    *miscellus*. Where both homeologs were present in at least ten *T. miscellus* Illumina reads,

286    and the observed allelic ratio was more balanced than 70:30 in either direction, we took

287    this as preliminary evidence that both homeologs were equally expressed.  In contrast,

288    any site where either the *T. pratensis* or *T.  dubius* parental homeolog was present at 10X

289    while the other was absent, was identified as suggestive of either complete silencing of

290    one parental homeolog or genomic homeolog loss.

291

292    **SNP validation**

293

294     A subset of SNPs identified using the above methods was analyzed using the Sequenom

295     MassARRAY iPLEX platform at the Center for Plant Genomics, Iowa State University.

296     Genomic DNA was extracted from leaf tissue of the three plants used for the

297     transcriptome sequencing, using a modified CTAB protocol (Doyle & Doyle 1987).

298     Multiplexed assays were designed using the Sequenom Assay Design 3.1 software for

299     four plexes containing a total of 139 SNPs between *T. dubius* and *T. pratensis*. Of these,

300     42 were scored as "potential gene loss" using the Illumina read data, 77 were scored as

301     "alleles balanced", 19 were scored as "low coverage in *T. miscellus*" and one had no *T.*

302     *miscellus* reads. This assay design was used to genotype a 384-well plate that included *T.*

303     *dubius*, *T. pratensis* and *T. miscellus* genomic DNA samples (~ 20 ng/ul). The resulting

304     data were analyzed using the MassARRAY Typer 4.0 Analyzer software. Using the

305     manufacturer's settings, the Sequenom software was used to call SNPs at "aggressive",

306     "moderate" and "conservative" degrees of confidence.

307

308 **Results**

309

310     **454 Sequencing, assembly and annotation of *T. dubius* cDNA sequences**

311

312     454 FLX sequencing of the normalized *T. dubius* cDNA pool from *T. dubius* leaf tissue

313     produced 822,594 reads (237 bp av. length) representing >195 MB of sequence. These

314     reads have been uploaded to the NCBI Short Read Archive (Accession Number

315    SRA009218.1). Assembly of the 454 FLX reads with the Roche 454 Newbler assembler

316    produced 33,515 contigs (14.7 Mb) with an average length of 439 bp (min = 96, Max =

317    3418), an average depth of 17.6 reads and N50 Contig Size of 626 bp (see Figure 1).

318

319    In comparison with other species in the Asteraceae, 21,498 (64 %) of the *T. dubius* 454

320    EST sequences matched previously characterized EST assemblies (TBLASTX) from

321    *Lactuca sativa, L. serriola* and *Helianthus annus* with p-values of e-5 or better. This low

322    percentage may reflect the low depth and coverage in many of our 454 contigs (Figure 1)

323    or significant divergence among the species. Of the 21,498 hits, 18,526 (86 %) were to

324    unique EST assemblies in this curated database. The 14 % of non-unique contigs may be

325    due to paralogous sequences in *T. dubius*, or to non-overlapping assemblies of *T. dubius*

326    sequence from the same cDNA template, as the 'shotgun' nature of 454 sequencing

327    enables simultaneous sampling of discrete template regions. The majority of best matches

328    occurred between *T. dubius* and *L. sativa* (Table 1). In comparison with *A. thaliana*,

329    18,923 *T. dubius* 454 EST contig assemblies match *A. thaliana* CDS sequences

330    (TBLASTX) at p-values of e-5 or better, while a total of 22,946 *T. dubius* contigs hit at

331    least one sequence in either the *A. thaliana* or the Asteraceae collection.

332

333    **SNP discovery**

334

335    Non-normalized cDNA pools sequenced on single lanes of an Illumina GAII analyzer

336    resulted in 7,128,226, 6,840,425 and 6,729,215 reads from *T. dubius, T. pratensis* and *T.*

337    *miscellus*, respectively. These reads have been uploaded to the NCBI Short Read Archive

338     (Accession Number SRA009218.1). Alignment of pooled Illumina reads to the *T. dubius*

339     454 assembled EST reference sequences with a mismatch tolerance of 2 bp followed by

340     identification of polymorphic sites that were represented to a minimum of three-fold

341     redundancy in both *T. dubius* and *T. pratensis* revealed >45,000 potential SNPs within

342     10,428 contigs.  To reduce the risk of misaligning repetitive or highly paralogous

343     sequences, parameters were adjusted to permit only a single mismatch over the length of

344     the Illumina reads. Of the total pooled *T. dubius*, *T. pratensis* and *T. miscellus* Illumina

345     reads, 11,050,022 (53.4%) aligned.  The remaining reads were unaligned because they

346     did not map a unique location in the 454 contig reference sequence collection, or they did

347     not meet the single mismatch criterion.  This higher confidence alignment, when parsed

348     for polymorphic sites that were represented to a minimum of three-fold redundancy in

349     both *T. dubius* and *T. pratensis*, resulted in the identification of 24,078 potential SNPs

350     between *T. dubius* and *T. pratensis* within 7,837 unique 454 EST contig reference

351     sequences. To identify an even higher-quality collection of potential SNP sites between

352     *T. dubius* and *T. pratensis*, the aforementioned alignments were parsed for SNP sites that

353     were represented to a minimum depth of 10X in both the *T. dubius* and *T. pratensis* data

354     sets. This high-quality collection that maximizes the likelihood that discovered

355     polymorphic sites represent true SNPs between *T. dubius* and *T. pratensis* consists of

356     7,782 SNPs within 2,885 unique contigs.

357

358     Of the 7,782 SNPs, 2,989 had sufficient *T. miscellus* Illumina reads for transcriptome

359     analysis. Of these, 2,064 (69 %) appeared to show equal homeolog expression in *T.*

360     *miscellus*, 671 (22 %) showed differential expression in *T. miscellus*, and 254 (8.5 %)

361    showed potential homeolog loss in *T. miscellus*. Interestingly, the cases of differential

362    expression were mainly due to higher expression of the *T. dubius* homeolog than of the *T.*

363    *pratensis* homeolog (454/671; 77%) and the majority of the apparent losses were also of

364    the *T. dubius* homeolog (164/254; 65%) rather than the *T. pratensis* homeolog (90/254).

365

366    **SNP validation**

367

368    Sequenom MassARRAY iPLEX  assays were designed for 139 of the putative SNPs

369    (four plexes). These assays were used to analyze the genomic DNA of the two diploid

370    plants whose transcriptomes were used for 454 and Illumina sequencing. For 19 of the

371    assays, the Sequenom assay failed to call a SNP in both diploid species, and 22 assays

372    only worked in one of the diploid species. This failure rate is comparable to those

373    obtained by other groups (Dunstan *et al.* 2007). Of the 98 informative assays (Table 2),

374    92 (94%) confirmed the SNP calls. In five of the remaining assays, the correct

375    polymorphism was present but there was an extra allele in the genome of one diploid (i.e.

376    heterozygosity) that had not been detected by via transcriptome sequencing. In only one

377    case did the base call differ between the sequencing and Sequenom methods: here

378    Sequenom indicated the same base in both alleles.

379

380    We then examined the Sequenom data for the genomic DNA of the *T. miscellus* plant. Of

381    the 139 SNP assays, 41 did not successfully call any bases within our confidence limits in

382    this plant. In 28 of these 41 cases, the assay also failed to call a SNP in one or both

383    diploid species, but in the remaining 13 cases, the assay called a SNP in both diploid

384   species but not in *T. miscellus* (see Table 2). In another 13 cases, one or more SNPs were

385   called in *T. miscellus*, but a base had only been successfully called in one of the diploids

386   (not shown in Table 2). Thus, in total, 85 of the 139 Sequenom assays (61 %) provided a

387   call. In no cases did we find a SNP homeolog present in *T. miscellus* that had not been

388   found in either *T. dubius* or *T. pratensis* at that locus.

389

390   Only the Sequenom data for 81 assays were used to infer homeolog loss in *T. miscellus*.

391   Of the 85 assays that worked in all three plants, three were excluded due to

392   heterozygosity in *T. dubius* and a fourth because of an identical call in both diploids. Of

393   the 81 assays used, 47 gave evidence in *T. miscellus* of both *T. dubius* and *T. pratensis*

394   SNP homeologs, and 34 gave evidence of only one SNP homeolog. Thus, 41% of the

395   SNP loci give evidence for homeolog loss. Of these, 25 (74%) showed loss of the *T.*

396   *dubius* homeolog, and nine (26%) showed loss of the *T. pratensis* homeolog. If we

397   increase stringency by omitting "aggressive" calls (i.e. less confident Sequenom calls),

398   we find that 69 assays gave a call; of these, 44 gave evidence of both SNP homeologs in

399   *T. miscellus*, and 26 gave evidence of only one SNP homeolog.

400

401   We then compared the Sequenom and Illumina sequence data for the *T. miscellus* plant,

402   to discover how often Illumina expression data had successfully identified a candidate for

403   genomic loss (Table 3). Illumina read counts were correct in 89 % of the cases where

404   there was depth of coverage above 10X per SNP homeolog, and the Sequenom calls were

405   at conservative, moderate and aggressive levels of confidence (listed in descending

406   order). Where Illumina read data had predicted "potential gene loss", this was shown by

407    Sequenom analysis in 23 of 27 cases (85%). In four cases, homeologs were detected in

408    the genomic DNA by Sequenom but not in the transcriptome by Illumina (i.e. they were

409    scored as "potential gene loss"). This may be due to homeolog silencing. In contrast, four

410    SNP homeologs were detected in the transcriptome by Illumina (i.e. they were scored

411    "alleles balanced") but were not found in the genome by Sequenom. Manual examination

412    of the mass spectrometer traces for these calls suggested that three of them, which had all

413    been called at the "aggressive" (lowest) level of confidence, did in fact have both

414    homeologs present in the gDNA. In no cases did a contradiction occur where Illumina

415    showed no expression of one homeolog and Sequenom loss of the other homeolog.

416

# Discussion

418

419    Genomic resources are scarce for many organisms that are studied in a natural ecological

420    or evolutionary context (Ellegren 2008; Hudson 2008). Here, we demonstrate a protocol

421    that uses next-generation technologies to rapidly develop SNP markers in many hundreds

422    of genes in a species which is a good evolutionary model but which until now has not

423    been a genetic model organism. These SNP markers have many potential uses. We have

424    used them to distinguish between homeologous genes in the recent natural allopolyploid

425    *T. miscellus*. Using transcriptome profiling and Sequenom genotyping, we have detected

426    many cases of gene loss. Below we discuss the biological implications of our findings in

427    *T. miscellus*, and the general utility of the methods described in this paper.

428

**Biological implications of findings in *T. miscellus***

430

431 This paper provides the first large-scale analysis of homeologous gene loss in a recent

432 (~40-generation-old) natural allopolyploid. In a single *T. miscellus* individual we found

433 254 cases of putative homeolog loss or silencing by transcriptome profiling with Illumina

434 sequencing (3% of all SNPs). Sequenom analysis confirmed that in a sample of 27 of

435 these SNPs, 23 (85%) were cases of genomic homeolog loss. The remaining 15% are

436 likely to be homeologs that are present in the genome but were not being expressed at the

437 time of sampling in the leaf tissue subject to transcriptome analysis. Homeolog loss

438 therefore appears to be more common than homeolog silencing (i.e. lack of expression of

439 a gene found in the genome) in this species.

440

441 We found preferential loss of *T. dubius* homeologs over *T. pratensis* homeologs in the

442 allopolyploid *T. miscellus* in this study. Illumina read data on the transcriptome suggested

443 loss or silencing of the *T. dubius* homeolog in 164 of 254 SNPs (64%) showing homeolog

444 loss or silencing, and Sequenom analysis of the genome suggested loss of the *T. dubius*

445 homeolog in 25 of 34 SNPs (73%) showing homeolog loss. In earlier studies a similar

446 bias was found: combined results from Buggs *et al*. (2009), Tate *et al*. (2006), and Tate *et*

447 *al*. (2009a) gave 56 *T. dubius* homeolog losses, and 27 *T. pratensis* homeolog losses

448 across multiple populations. Interestingly, we also found a bias in gene expression in our

449 Illumina read data, with *T. dubius* homeologs tending to be expressed more than *T.*

450 *pratensis* homeologs in 77% of the SNPs where we detected differential expression.

451     Because *T. dubius* ESTs were used as the reference sequence we might expect a bias

452     towards the alignment of Illumina reads derived from *T. dubius* homeologs. This possible

453     bias may have contributed to the apparent higher expression of the *T. dubius* homeolog at

454     many loci, but also suggests that the finding of a higher rate of loss of *T. dubius*

455     homeologs is a robust result.

456

457     It is notable that a similar bias towards loss of *T. dubius* genetic material and higher

458     expression of *T. dubius* genes has been found for rDNA in both *T. miscellus* and *T. mirus*,

459     an allopolyploid that has *T. dubius* as the paternal parent and *T. porrifolius* as the

460     maternal parent (Kovarik *et al.* 2005), In both species, concerted evolution has reduced

461     the copy numbers of rDNA units derived mainly from the *T. dubius* diploid parent but,

462     paradoxically, repeats of *T. dubius* origin dominate transcription in most populations

463     studied (Matyasek *et al.* 2007). *Tragopogon mirus* also shows a bias toward loss of *T.*

464     *dubius* homeologs using CAPS markers (Koh *et al*., submitted)

465

466     What causes the bias towards higher rates of gene loss and increased expression of *T.*

467     *dubius* homeologs? One possibility might be maternal effects due to cytoplasmic-nuclear

468     interactions. The *T. miscellus* plant in the current study, as well as all *T. mirus* plants and

469     the majority of *T. miscellus* plants included in other studies, has *T. dubius* as the paternal

470     parent. Perhaps selection favors maintaining ancestral similarity in the cytoplasmic and

471     nuclear genomes. Another explanation might be the higher genetic variability of *T.*

472     *dubius* populations (Soltis *et al.* 1995); it is possible that the *T. dubius* individual that we

473     examined from Oakesdale was not genetically identical to the actual *T. dubius* progenitor

474    of *T. miscellus* from Oakesdale. However, it seems unlikely that the bias is due to the

475    selection of an inappropriate *T. dubius* genotype in this study as the other studies cited

476    above as showing the same pattern have examined multiple *T. dubius* individuals. Our

477    results also agree with those found in other species. In synthetic allopolyploids of

478    *Brassica*, genomic changes occur more often in the paternal genome (Song *et al.* 1995).

479    In natural *Gossypium hirsutum* (Flagel *et al.* 2008) and synthetic *Arabidopsis*

480    allopolyploids (Wang *et al.* 2006), homeolog expression biases also tend to be in favor of

481    the paternal genome.  In maize, it has recently been shown that paternal genomic

482    imprinting influences gene expression patterns in hybrids (Swanson-Wagner *et al.* 2009).

483

484    One mechanism by which homeolog loss may occur in *T. miscellus* is homeologous

485    recombination, in which fragments of chromosomes can be lost. Ownbey (1950)

486    observed multivalent formation in early generations of natural *T. miscellus*, and rare

487    patterns of isozyme variation in *T. miscellus* are consistent with homeologous

488    recombination (Soltis *et al.* 1995). More recently, Lim *et al*. (2008) and Tate *et al.*

489    (2009b) report multivalent formation in both natural and synthetic *Tragopogon*

490    allopolyploids, along with unisomy, trisomy, and reciprocal translocations in natural

491    *Tragopogon* allopolyploids. Homoeologous recombination appears to have caused loss of

492    chromosome fragments in re-synthesised *Brassica* allopolyploids (Song *et al.* 1995;

493    Gaeta *et al.* 2007). Another possible mechanism of homeolog loss is gene conversion, as

494    has been found for rRNA genes in both *T. miscellus* and *T. mirus* (Kovarik *et al.* 2005;

495    Matyasek *et al.* 2007).

496

497    High-throughput SNP discovery together with the genotyping of many natural *T.*

498    *miscellus* plants of independent origin and $F_1$ hybrids will enable us to examine genome-

499    wide patterns of homeolog loss in this species. As SNPs are abundant in many species

500    and easily detected (Gut 2001; Kwok 2001), they are excellent genetic markers for the

501    generation of dense genetic maps that can support  marker-assisted selection (MAS) and

502    association genetics programs, as well as inform on genome organization and function

503    (Pavy *et al.* 2008; Slate *et al.* 2009). In *T. miscellus*, application of these markers will

504    enable us to understand further the causes of homeolog loss in this allopolyploid,

505    showing us whether or not homeolog losses occur in linkage groups – implying the loss

506    of large fragments of chromosomes – or in small fragments scattered throughout the

507    genome.

508

509    **Utility of methods**

510

511    In the space of a few months, we have been able to identify at high stringency 7,782

512    homeolog-specific SNP markers within 2,885 unique contigs in *T. miscellus* using next-

513    generation sequencing. The number of homeologous genes available for study has

514    therefore been increased by two orders of magnitude compared to previous studies using

515    a "one gene at a time" approach (Tate *et al.* 2006; Buggs *et al.* 2009; Tate *et al.* 2009a).

516    The number of actual SNPs discovered is likely to be much higher than this, as we were

517    likely over-stringent. We have developed working assays for 85 of these SNPs using

518    Sequenom MassARRAY iPLEX technology. This high-throughput approach transforms

519    our ability to study molecular evolution in *T. miscellus*.

520

521    The use of transcriptome sequencing with polyA purification is valuable for targeting

522    functional genes for SNP discovery, as clearly shown by this study. However, there is the

523    possibility that when these markers are then used to study the genome, polymorphisms

524    will be discovered due to the presence of silent homeologs. In a few cases we found this:

525    six of 139 Sequenom SNP assays found polymorphisms in genomic DNA of diploid

526    plants that had not been detected by Illumina sequencing in the transcriptome. This was

527    an acceptably low level of polymorphism that was undiscovered by transcriptome

528    sequencing. However, it should be noted that *T. dubius* and *T. pratensis* are mostly

529    selfing species (Cook & Soltis 1999; Cook & Soltis 2000) with limited polymorphism in

530    their introduced ranges in North America (Soltis *et al.* 1995; Symonds *et al.* 2009).

531    Outcrossing species with high heterozygosity may pose more difficulties in analysis.

532

533    Sequenom MassARRAY Typer 4.0 Analyzer software uses a three-parameter model to

534    calculate the significance of each putative genotype. This compares the size of peaks for

535    the possible bases at each SNP site and the peak for the unextended primer. Where an

536    assay is not working well, the non-extended primer will be found in greater abundance

537    than the extended oligonucleotides. For genotypes which are called, the degree of

538    confidence that can be placed on the call is described as "conservative", "moderate", or

539    "aggressive" in the software output. We found that four calls (three called at the

540    "aggressive" (lowest) level of confidence and one at the "moderate" level) were not

541    reliable due to failure to detect a base that was in fact present (a false negative). Manual

542    examination of the mass-spectrometer trace in most cases allowed the call to be

543    corrected.

544

545    This "false negative" problem is likely to be due to the malfunction of these specific

546    assays, rather than the reliability of "aggressive" calls in general. Certain assays can

547    function well in calling different bases in homozygotes, but in a heterozygote the primers

548    bind preferentially to one allele, resulting in a false homozygote call. One reason why this

549    occurs is if there is another SNP close to the SNP site that is being assayed (Liu *et al.*

550    2009). Preferential binding of primers can be assessed by genotyping more individuals

551    that are expected to be heterozygous. If they all appear to be homozygous, then the

552    Sequenom assay for that SNP should be rejected. We did this (see below) and found that

553    these assays did not work correctly in multiple individuals. In addition, if we discard all

554    aggressive Sequenom calls, we find that the correspondence between the Illumina and

555    Sequenom data rises only slightly from 89% to 93%. This also suggests that there is not a

556    general problem with the reliability of "aggressive" calls.

557

558    This study also demonstrates that transcriptome profiling using Illumina sequencing is a

559    useful method for identifying candidate homeologs for the study of homeolog loss in an

560    allopolyploid species. This allows us to target these genes for developing SNP-typing

561    assays, saving both time and money. The major cost in using Sequenom genotyping is the

562    production of primers. Each SNP requires three primers: two for an initial amplification

563    of the target region and one for the SNP-typing reaction. Once these primers have been

564    synthesized, many samples can be SNP-typed at relatively low cost. We made use of this

565     fact by screening an additional 94 individuals: a total of 87 diploid and *T. miscellus* plants

566     from five natural populations, two 50-year-old herbarium specimens and five artificial

567     crosses. Preliminary analyses of this survey allowed us to identify polymorphisms in the

568     diploid plants and calculate allelic diversity. This data set showed repeatability of some

569     homeolog losses in natural *T. miscellus* populations of different origins. Finally, we also

570     found the first evidence for rare loss of alleles in $F_1$ hybrids between *T. dubius* and *T.*

571     *pratensis*. Robust analysis of this data set is ongoing.

572

573     **Broader applicability**

574

575     Transcriptome sequencing by 454 has many potential applications in ecology (Ellegren

576     2008; Wang *et al.* 2009; Wheat 2008). It has been used for the *de novo* characterization

577     of the transcriptome of the Glanville fritillary butterfly (Vera *et al.* 2008) and the

578     *Eucalyptus grandis* genome (Novaes *et al.* 2008). Recent work in model organisms has

579     used short-read sequencing to study differences in expression of SNP-containing alleles,

580     for example in micro-RNAs in mice (Kim & Bartel 2009). Sequenom MassARRAY

581     genotyping has been used to study allelic expression in hybrid maize (Stupar & Springer

582     2006) and levels of homeolog expression in allopolyploid cotton (Flagel *et al.* 2008;

583     Chaudhary *et al.* 2009; Flagel *et al.* 2009). This study demonstrates the effectiveness of a

584     hybrid Illumina and 454 sequencing approach and Sequenom MassARRAY iPLEX

585     genotyping to increase dramatically our ability to study the evolution of duplicated genes

586     in natural allopolyploids such as *T. miscellus.* These methods could be applied to any

587     organism, allowing efficient and cost-effective generation of SNP markers.

588

589

590

# Acknowledgments

596

# References

Adams KL, Cronn R, Percifield R, Wendel JF (2003) Genes duplicated by polyploidy
show unequal contributions to the transcriptome and organ-specific reciprocal
silencing. *Proceedings of the National Academy of Sciences of the United States
of America*, **100**, 4649-4654.

Adams KL, Wendel JF (2004) Exploring the genomic mysteries of polyploidy in cotton.
*Biological Journal of the Linnean Society*, **82**, 573-581.

Adams KL, Wendel JF (2005) Polyploidy and genome evolution in plants. *Current
Opinion in Plant Biology*, **8**, 135-141.

Ainouche ML, Baumel A, Salmon A (2004) *Spartina anglica* C. E. Hubbard: a natural
model system for analysing early evolutionary changes that affect allopolyploid
genomes. *Biological Journal of the Linnean Society*, **82**, 475-484.

609    Barbazuk WB, Emrich SJ, Chen HD, Li L, Schnable PS (2007) SNP discovery via 454

610        transcriptome sequencing. *Plant Journal*, **51**, 910-918.

611    Bottley A, Xia GM, Koebner RMD (2006) Homoeologous gene silencing in hexaploid

612        wheat. *Plant Journal*, **47**, 897-906.

613    Buggs RJA (2008) Towards natural polyploid model organisms. *Molecular Ecology*, **17**,

614        1875-1876.

615    Buggs RJA, Doust AN, Tate JA*, et al.* (2009) Gene loss and silencing in *Tragopogon*

616        *miscellus* (Asteraceae): comparison of natural and synthetic allotetraploids.

617        *Heredity*, **103**, 73-81.

618    Buggs RJA, Pannell JR (2007) Ecological differentiation and diploid superiority across a

619        moving ploidy contact zone. *Evolution*, **61**, 125 - 140.

620    Chaisson MJ, Pevzner PA (2008) Short read fragment assembly of bacterial genomes.

621        *Genome Research*, **18**, 324-330.

622    Chaudhary B, Flagel L, Stupar RM*, et al.* (2009) Reciprocal silencing, transcriptional

623        bias and functional divergence of homeologs in polyploid cotton (*Gossypium*).

624        *Genetics*, **182**, 503-517.

625    Chen M, Ha M, Lackey E, Wang JL, Chen ZJ (2008) RNAi of met1 reduces DNA

626        methylation and induces genome-specific changes in gene expression and

627        centromeric small RNA accumulation in *Arabidopsis* allopolyploids. *Genetics*,

628        **178**, 1845-1858.

629    Chen ZJ, Wang J, Tian L*, et al.* (2004) The development of an *Arabidopsis* model system

630        for genome-wide analysis of polyploidy effects. *Biological Journal of the Linnean*

631        *Society*, **82**, 689-700.

632    Cheung F, Haas B, Goldberg S*, et al.* (2006) Sequencing *Medicago truncatula* expressed

633            sequenced tags using 454 Life Sciences technology. *BMC Genomics*, **7**, 272.

634    Cheung F, Win J, Lang J*, et al.* (2008) Analysis of the *Pythium ultimum* transcriptome

635            using Sanger and Pyrosequencing approaches. *BMC Genomics*, **9**, 542.

636    Cook LM, Soltis PS (1999) Mating systems of diploid and allotetraploid populations of

637            *Tragopogon* (Asteraceae). I. Natural populations. *Heredity*, **82**, 237-244.

638    Cook LM, Soltis PS (2000) Mating systems of diploid and allotetraploid populations of

639            Tragopogon (Asteraceae). II. Artificial populations. *Heredity*, **84**, 410-415.

640    Dong Y, Liu Z, Shan X*, et al.* (2005) Allopolyploidy in wheat induces rapid and heritable

641            alterations in DNA methylation patterns of cellular genes and mobile elements.

642            *Russian Journal of Genetics*, **41**, 890-896.

643    Doyle J, Doyle JL (1987) Genomic plant DNA preparation from fresh tissue-CTAB

644            method. *Phytochemical Bulletin*, **19**, 11-15.

645    Dunstan S, Hue N, Rockett K*, et al.* (2007) A TNF region haplotype offers protection

646            from typhoid fever in Vietnamese patients. *Human Genetics*, **122**, 51-61.

647    Ellegren H (2008) Sequencing goes 454 and takes large-scale genomics into the wild.

648            *Molecular Ecology*, **17**, 1629-1631.

649    Emrich SJ, Barbazuk WB, Li L, Schnable PS (2007) Gene discovery and annotation

650            using LCM-454 transcriptome sequencing. *Genome Research*, **17**, 69-73.

651    Eveland AL, McCarty DR, Koch KE (2008) Transcript profiling by 3'-untranslated

652            region sequencing resolves expression of gene families. *Plant Physiology*, **146**,

653            32-44.

654 Feldman M, Liu B, Segal G, *et al.* (1997) Rapid elimination of low-copy DNA sequences

655         in polyploid wheat: A possible mechanism for differentiation of homoeologous

656         chromosomes. *Genetics*, **147**, 1381-1387.

657 Flagel LE, Chen L, Chaudhary B, Wendel JF (2009) Coordinated and fine-scale control

658         of homoeologous gene expression in allotetraploid cotton. *Journal of Heredity*,

659         **100**, 487-490.

660 Flagel LE, Udall J, Nettleton D, Wendel J (2008) Duplicate gene expression in

661         allopolyploid *Gossypium* reveals two temporally distinct phases of expression

662         evolution. *BMC Biology*, **6**, 16.

663 Gabriel S, Ziaugra L, Tabbaa D (2009) SNP genotyping using the Sequenom

664         MassARRAY iPLEX platform. *Current Protocols in Human Genetics*, **60**,

665         2.12.11-12.12.18.

666 Gaeta RT, Pires JC, Iniguez-Luy F, Leon E, Osborn TC (2007) Genomic changes in

667         resynthesized *Brassica napus* and their effect on gene expression and phenotype.

668         *Plant Cell*, **19**, 3403-3417

669 Guo M, Rupe MA, Zinselmeier C, *et al.* (2004) Allelic variation of gene expression in

670         maize hybrids. *Plant Cell*, **16**, 1707-1716.

671 Gut IG (2001) Automation in genotyping of single nucleotide polymorphisms. *Human*

672         *Mutation*, **17**, 475-492.

673 Hegarty M, Barker G, Wilson I, *et al.* (2006) Transcriptome shock after interspecific

674         hybridization in *Senecio* is ameliorated by genome duplication. *Current Biology*,

675         **16**, 1652-1659.

676    Hegarty M, Jones J, Wilson I, *et al.* (2005) Development of anonymous cDNA

677        microarrays to study changes to the *Senecio* floral transcriptome during hybrid

678        speciation. *Molecular Ecology*, **14**, 2493-2510.

679    Hillier LW, Marth GT, Quinlan AR, *et al.* (2008) Whole-genome sequencing and variant

680        discovery in *C. elegans*. *Nature Methods*, **5**, 183-188.

681    Holt RA, Jones SJM (2008) The new paradigm of flow cell sequencing. *Genome*

682        *Research*, **18**, 839-846.

683    Hudson ME (2008) Sequencing breakthroughs for genomic ecology and evolutionary

684        biology. *Molecular Ecology Resources*, **8**, 3-17.

685    Johnson DS, Mortazavi A, Myers RM, Wold B (2007) Genome-wide mapping of in vivo

686        protein-DNA interactions. *Science*, **316**, 1497-1502.

687    Joly S, Rauscher JT, Sherman-Broyles SL, Brown AHD, Doyle JJ (2004) Evolutionary

688        dynamics and preferential expression of homeologous 18S-5.8S-26S nuclear

689        ribosomal genes in natural and artificial Glycine allopolyploids. *Molecular*

690        *Biology and Evolution*, **21**, 1409-1421.

691    Judd WS, Campbell CS, Kellogg EA, Stevens PF, Donoghue MJ (2007) *Plant*

692        *Systematics: A Phylogenetic Approach*, Third Edition edn. Sinauer, Sunderland,

693        MA.

694    Kashkush K, Feldman M, Levy AA (2002) Gene loss, silencing and activation in a newly

695        synthesized wheat allotetraploid. *Genetics*, **160**, 1651-1659.

696    Kihara H, Ono T (1927) Chromosomenzahlen und systematische Gruppierung der

697        *Rumex*-Arten. *Zeitschrift für Zellforschung und mikroskopische Anatomie*, **4**, 475-

698        481.

699     Kim J, Bartel DP (2009) Allelic imbalance sequencing reveals that single-nucleotide

700         polymorphisms frequently alter microRNA-directed repression. *Nature*

701         *Biotechnology*, **27**, 472-477.

702     Kim S, Yoo M-J, Albert VA*, et al.* (2004) Phylogeny and diversification of B-function

703         MADS-box genes in angiosperms: evolutionary and functional implications of a

704         260-million-year-old duplication. *American Journal of Botany*, **91**, 2102-2118.

705     Kovarik A, Pires JC, Leitch AR*, et al.* (2005) Rapid concerted evolution of nuclear

706         ribosomal DNA in two *Tragopogon* allopolyploids of recent and recurrent origin.

707         *Genetics*, **169**, 931-944.

708     Kwok PY (2001) Methods for genotyping single nucleotide polymorphisms. *Annual*

709         *Review of Genomics and Human Genetics*, **2**, 235-258.

710     Levy AA, Feldman M (2004) Genetic and epigenetic reprogramming of the wheat

711         genome upon allopolyploidization. *Biological Journal of the Linnean Society*, **82**,

712         607-613.

713     Lim KY, Matyasek R, Kovarik A, Leitch AR (2004) Genome evolution in allotetraploid

714         *Nicotiana*. *Biological Journal of the Linnean Society*, **82**, 599-606.

715     Lim KY, Soltis DE, Soltis PS*, et al.* (2008) Rapid chromosome evolution in recently

716         formed polyploids in *Tragopogon* (Asteraceae). *PLoS ONE*, **3**, e3353.

717     Liu S, Chen D, Makarevitch I*, et al.* (2009) High-throughput genetic mapping of mutants

718         via quantitative SNP-typing. *Genetics*, **submitted**.

719     Lynch M, Conery JS (2000) The evolutionary fate and consequences of duplicate genes.

720         *Science*, **290**, 1151-1155.

721    Marth GT, Korf I, Yandell MD*, et al.* (1999) A general approach to single-nucleotide

722         polymorphism discovery. *Nat Genet*, **23**, 452-456.

723    Matyasek R, Tate JA, Lim YK*, et al.* (2007) Concerted evolution of rDNA in recently

724         formed *Tragopogon* allotetraploids is typically associated with an inverse

725         correlation between gene copy number and expression. *Genetics*, **176**, 2509-2519.

726    Mavrodiev EV, Tancig M, Sherwood AM*, et al.* (2005) Phylogeny of *Tragopogon* L.

727         (Asteraceae) based on internal and external transcribed spacer sequence data.

728         *International Journal of Plant Sciences*, **166**, 117-133.

729    Novaes E, Drost D, Farmerie W*, et al.* (2008) High-throughput gene and SNP discovery

730         in *Eucalyptus grandis*, an uncharacterized genome. *BMC Genomics*, **9**, 312.

731    Novak SJ, Soltis DE, Soltis PS (1991) Ownbey *Tragopogons* - 40 Years Later. *American*

732         *Journal of Botany*, **78**, 1586-1600.

733    Ownbey M (1950) Natural hybridization and amphiploidy in the genus *Tragopogon*.

734         *American Journal of Botany*, **37**, 487-499.

735    Pavy N, Pelgas B, Beauseigle S*, et al.* (2008) Enhancing genetic mapping of complex

736         genomes through the design of highly-multiplexed SNP arrays: application to the

737         large and unsequenced genomes of white spruce and black spruce. *BMC*

738         *Genomics*, **9**, 21.

739    Petit M, Lim K, Julio E*, et al.* (2007) Differential impact of retrotransposon populations

740         on the genome of allotetraploid tobacco (*Nicotiana tabacum*). *Molecular Genetics*

741         *and Genomics*, **278**, 1-15.

742     Quackenbush J, Liang F, Holt I, Pertea G, Upton J (2000) The TIGR gene indices:

743             reconstruction and representation of expressed gene sequences. *Nucleic Acids*

744             *Res*, **28**, 141-145.

745     Sadava DE, Heller HC, Orians GH, Purves WK, Hillis DM (2008) *Life: the science of*

746             *biology*, Eighth edn. Macmillan, New York.

747     Slate J, Gratten J, Beraldi D*, et al.* (2009) Gene mapping in the wild with SNPs:

748             guidelines and future directions. *Genetica*, **136**, 97-107.

749     Soltis DE, Soltis PS (1989) Allopolyploid speciation in *Tragopogon*: insights from

750             chloroplast DNA. *American Journal of Botany*, **76**, 14-18.

751     Soltis DE, Soltis PS (1999) Polyploidy: recurrent formation and genome evolution.

752             *Trends in Ecology & Evolution*, **14**, 348-352.

753     Soltis DE, Soltis PS, Pires JC*, et al.* (2004a) Recent and recurrent polyploidy in

754             *Tragopogon* (Asteraceae): cytogenetic, genomic and genetic comparisons.

755             *Biological Journal of the Linnean Society*, **82**, 485-501.

756     Soltis DE, Soltis PS, Tate JA (2004b) Advances in the study of polyploidy since Plant

757             Speciation. *New Phytologist*, **161**, 173-191.

758     Soltis PS, Plunkett GM, Novak SJ, Soltis DE (1995) Genetic variation in *Tragopogon*

759             species: additional origins of the allotetraploids *T. mirus* and *T. miscellus*

760             (Compositae). *American Journal of Botany*, **82**, 1329-1341.

761     Song K, Lu P, Tang K, Osborn TC (1995) Rapid genome change in synthetic polyploids

762             of *Brassica* and its implications for polyploid evolution. *Proceedings of the*

763             *National Academy of Sciences of the United States of America*, **92**, 7719-7723.

764    Stebbins GL (1942) Polyploid Complexes in Relation to Ecology and the History of

765         Floras. *American Naturalist*, **76**, 36-45.

766    Stupar RM, Springer NM (2006) *Cis*-transcriptional variation in maize inbred lines B73

767         and Mo17 leads to additive expression patterns in the F-1 hybrid. *Genetics*, **173**,

768         2199-2210.

769    Swanson-Wagner R, DeCook R, Jia Y*, et al.* (2009) Widespread paternal genomic

770         imprinting of trans-eQTL influences gene expression patterns in maize hybrids.

771         *Science*, **submitted**

772    Symonds VV, Soltis PS, Soltis DE (2009) The dynamics of polyploid formation in

773         *Tragopogon* (Asteraceae): recurrent formation, gene flow, and population

774         structure. *Evolution*, **submitted**.

775    Tate J, Joshi P, Soltis K, Soltis P, Soltis D (2009a) On the road to diploidization?

776         Homoeolog loss in independently formed populations of the allopolyploid

777         *Tragopogon miscellus* (Asteraceae). *BMC Plant Biology*, **9**, 80.

778    Tate JA, Ni ZF, Scheen AC*, et al.* (2006) Evolution and expression of homeologous loci

779         in *Tragopogon miscellus* (Asteraceae), a recent and reciprocally formed

780         allopolyploid. *Genetics*, **173**, 1599-1611.

781    Tate JA, Symonds VV, Doust AN*, et al.* (2009b) Synthetic polyploids of *Tragopogon*

782         *miscellus* and *T. mirus* (Asteraceae): 60 Years after Ownbey's discovery. *Am. J.*

783         *Bot.*, **96**, 979-988.

784    Udall JA, Wendel JF (2006) Polyploidy and crop improvement. *Crop Science*, **46**, S3-

785         S14.

786     Van Tassell CP, Smith TPL, Matukumalli LK*, et al.* (2008) SNP discovery and allele

787         frequency estimation by deep sequencing of reduced representation libraries.

788         *Nature Methods*, **5**, 247-252.

789     Vera JC, Wheat CW, Fescemyer HW*, et al.* (2008) Rapid transcriptome characterization

790         for a nonmodel organism using 454 pyrosequencing. *Molecular Ecology*, **17**,

791         1636-1647.

792     Wang J, Tian L, Lee H-S*, et al.* (2006) Genomewide nonadditive gene regulation in

793         *Arabidopsis* allotetraploids. *Genetics*, **172**, 507-517.

794     Wang Z, Gerstein M, Snyder M (2009) RNA-Seq: a revolutionary tool for

795         transcriptomics. *Nature Reviews Genetics*, **10**, 57-63.

796     Wheat C (2008) Rapidly developing functional genomics in ecological model systems via

797         454 transcriptome sequencing. *Genetica*, **Online First 18 October 2008**

798
799
800

800 **Table 1: Results of *Tragopogon dubius* similarity searches (BLASTX).**
801

| Sequence collection used for similarity searches | Number of 454 contigs with similarity @ 1e-05 | Number of Annotation sequences hit @ 1e-05 | Number of 454 contigs with similarity @ 1e-10 | Number of Annotation sequences hit @ 1e-10 |
|---|---|---|---|---|
| *Lactuca sativa, Lactuca serriola and Helianthus annus* Gene Index | 21,498 (*Lactuca sativa*: 11,080 *Lactuca serriola*: 6078 *Helianthus annus*: 4340) | 16,611 | 18,526 (*Lactuca sativa*: 9,731 *Lactuca serriola*: 5264 *Helianthus annus*: 3531) | 14,914 |
| Arabidopsis annotated peptides | 18,923 | 11,086 | 16,412 | 10,180 |

802