

1 A draft genome of *Alliaria petiolata* (garlic
2 mustard) as a model system for invasion
3 genetics
4

5 Nikolay Alabi¹, Yihan Wu¹, Oliver Bossdorf², Loren H. Rieseberg³
6 and Robert I. Colautti^{1*}
7

8 ¹BIOLOGY DEPARTMENT, QUEEN'S UNIVERSITY, KINGSTON, ONTARIO, CANADA

9 ²INSTITUTE OF ECOLOGY AND EVOLUTION, UNIVERSITY OF TÜBINGEN, TÜBINGEN, GERMANY

10 ³DEPARTMENT OF BOTANY AND BIODIVERSITY RESEARCH CENTRE, UNIVERSITY OF BRITISH COLUMBIA, VANCOUVER, BC,
11 CANADA

12

13 CORRESPONDING AUTHOR: ROBERT.COLAUTTI@QUEENSU.CA

14

15

16

17 **Keywords:** *Alliaria petiolata*, EFCC3, invasion genetics, garlic mustard, Illumina, mate pairs

18 Abstract

19 The emerging field of invasion genetics examines the genetic causes and consequences of biological
20 invasions, but few study systems are available that integrate deep ecological knowledge with genomic
21 tools. Here we report on the *de novo* assembly and annotation of a genome for the biennial herb *Alliaria*
22 *petiolata* (M. Bieb.) Cavara & Grande (Brassicaceae), which is widespread in Eurasia and invasive across
23 much of temperate North America. Our goal was to sequence and annotate a genome to complement
24 resources available from hundreds of published ecological studies, a global field survey, and hundreds of
25 genetic lines maintained in Germany and Canada. We sequenced a genotype (EFCC-3-20) collected from
26 the native range near Venice, Italy and sequenced paired-end and mate pair libraries at ~70× coverage.
27 A *de novo* assembly resulted in a highly continuous draft genome ($N_{50} = 121\text{Mb}$; $L_{50} = 2$) with 99.7% of
28 the 1.1Gb genome mapping to contigs of at least 50Kb in length. A total of 64,770 predicted genes in the
29 annotated genome include 99% of plant BUSCO genes and 98% of transcriptome reads. Consistent with
30 previous reports of (auto)hexaploidy in western Europe. Almost, we found that almost one third of
31 BUSCO genes (390/1440) mapped to two or more scaffolds despite a genome-wide average of < 2%
32 heterozygosity. The continuity and gene space quality of our draft genome assembly will enable
33 genomic studies of *A. petiolata* to address questions relevant to invasion genetics and conservation
34 efforts.

35 Introduction

36 Biological invasions are a threat to global biodiversity with significant impacts to human health and
37 welfare (Mack et al., 2000). They also present opportunities for large-scale ‘natural’ experiments to
38 study ecological and evolutionary processes in the wild (Mooney and Cleland, 2001; Sax et al., 2005).
39 Despite a large body of ecological research and a growing number of evolutionary studies of invasive
40 species, functional genetic and ‘omics approaches have been rare among studies of invasive species
41 until recently (reviewed in (Barrett, 2015; Bock et al., 2015); but see e.g. (Barrett et al., 2016; Boheemen
42 et al., 2017; Bourne et al., 2020)). Apart from studies using neutral markers to assess population
43 structure (reviewed in (Dlugosch and Parker, 2008)), relatively little is known about the genetic causes
44 and consequences of biotic invasions at the molecular level. A lack of genomic resources has hindered
45 high-resolution genetic studies of most invasive species. Here, we report on a draft genome for the
46 herbaceous biennial plant *Alliaria petiolata*, a plant invader in North America with potential to become a
47 model system for the emerging field of invasion genetics.

48 Several factors favour *A. petiolata* as an emerging model system for invasion genetics (Colautti et al.,
49 2014). First, it is widely distributed with variable ecological impacts throughout its range (Lankau et al.,
50 2009; USDA, 2020). Second, it has a relatively simple two-year life cycle with well-defined life stages
51 (seed, rosette, bolting, senescent). This simple life history facilitates measurements of lifetime fitness,
52 natural selection, and their impacts on population dynamics in natural populations. Third, *A. petiolata* is
53 a member of the Brassicaceae and therefore benefits from genetic resources available for well-studied
54 species like *Brassica rapa* (canola) and the model plant *Arabidopsis thaliana*, providing opportunities for
55 functional and comparative genomics. Fourth, high selfing rates produce naturally inbred seed families
56 that can be maintained through single-seed descent. Fifth, *A. petiolata* has been the focal species in
57 hundreds of field surveys and experimental studies, including influential studies testing the role of
58 natural enemies (Lewis et al., 2006), competition (Prati and Bossdorf, 2004), the ‘novel weapons’
59 hypothesis (Callaway et al., 2008), competitive ability (Bossdorf et al., 2004), glucosinolate metabolism

60 (Haribal et al., 2001), and eco-evolutionary dynamics (Lankau et al., 2009). One mechanism that is
61 especially well-studied is the production of secondary metabolites and their effects on soil microbiota,
62 particularly the suppression of mycorrhizal fungi that form beneficial symbiotic networks among native
63 plant roots (Anthony et al., 2017; Duchesneau et al., 2020). Understanding the genomic basis of such
64 interactions with soil microbes will not only advance basic science but it also has potential applications
65 in plant restoration and agriculture.

66 Recent efforts to develop *A. petiolata* as a model system include the Global Garlic Mustard Field Survey
67 (GGMFS), which mobilized 164 participants from 16 countries across North America and Europe to
68 collect field data and seed samples across Europe and North America, resulting in thousands of seed
69 families from 383 distinct populations (Colautti et al., 2014). A subset of inbred lines collected across
70 North America have been maintained through single-line descent in labs in Germany (Bossdorf Lab,
71 University of Tübingen) and Canada (Colautti Lab, Queen's University). Adding to these resources, we
72 here report on a draft genome of a single *A. petiolata* genotype from Europe, annotated with RNA
73 sequencing of leaf and root tissue.

74 Methods & Materials

75 Study Organism and Line Derivation

76 *Alliaria petiolata* (M. Bieb.) Cavara & Grande is a biennial herbaceous plant in the Thlaspideae tribe of
77 the Brassicaceae family. It was introduced to North America prior to 1868, when it was discovered in
78 Long Island, New York (Nuzzo, 1993). By 1948 it was reported on the West Coast of North America and
79 has established in at least 37 U.S. States and five Canadian provinces from the Atlantic coast to the
80 Pacific (Cavers et al., 1979; USDA, 2020). As the only species of *Alliaria* with a broad distribution, *A.*
81 *petiolata* is relatively easy to identify in natural habitats owing to its white flowers and dentated peltate
82 leaves with long petioles. It is considered a noxious weed across most of its introduced range, due in
83 part to impacts on native plants and tree regeneration in deciduous forest ecosystems (Cipollini and
84 Cipollini, 2016; Stinson et al., 2007). Two genome types have been identified, including diploids ($2n =$
85 14) in Eastern Europe and Western Asia and hexaploids ($2n = 42$ and $2n = 36$) in Central Europe and
86 North America (Esmailbegi et al., 2018; Weiss-Schneeweiss and Schneeweiss, 2003).

87 S_0 generation

88 All source material for genome and transcriptome sequencing originated from a single individual grown
89 from seed collected as part of the GGMFS (Colautti et al., 2014). The inbred line used in this study is
90 from population EFCC3, collected in 2011 from a small forest fragment ($\sim 2,000\text{m}^2$) surrounded by
91 agricultural land, about 75km northwest of Venice, Italy (UTM 45.71°N \times 11.72°E). The specific seed line
92 was sampled at 3.2m along a 10m sampling transect originating at the edge of population EFCC3. This
93 inbred line (code EFCC3-3-20) is currently maintained along with other GGMFS seed collections by two of
94 the coauthors in replicate collections in Tübingen, Germany (Bossdorf) and Kingston, Ontario, Canada
95 (Colautti).

96 S_1 generation

97 In July 2012, ten seeds of the EFCC3-3-20 genotype from the original S_0 field collection were removed
98 from cold storage (4 °C), surface washed with a mild detergent and rinsed with distilled H₂O before
99 surface sterilizing in 10% bleach for 10 minutes. Sterilized seeds were again rinsed with distilled H₂O
100 before placing on filter paper saturated with distilled H₂O and sealed in a petri dish with paraffin wax.

101 We stratified seeds in the dark at 10 °C for ~90 days and thereafter inspected weekly until emerging
102 radicles were observed. Germinating seeds were transplanted into 4" plastic pots containing a peat soil
103 mixed with vermiculite that was watered to saturation and placed under shade cloth in the Horticulture
104 Greenhouses at the University of British Columbia. We let seedlings establish in soil, watered as needed,
105 for four weeks before a small amount (~5mm x 2mm) of young leaf meristem tissue was harvested from
106 a single individual and immediately preserved in liquid nitrogen. Roughly 25 to 50mm³ of this tissue was
107 divided into two separate 2mL screw-cap tubes, each containing two stainless steel ball bearings of 2mm
108 diameter. Tubes were flash-frozen in liquid nitrogen and used for genomic DNA purification and genome
109 sequencing.

110 A second individual from the same inbred family was transplanted to an 8" plastic pot and fertilized with
111 20/20/20 N/P/K fertilizer to encourage rosette growth before being moved outside from October 2012
112 to April 2013 for cold vernalization at the University of British Columbia Horticulture Greenhouses. In
113 April 2013, we moved the plant back into the greenhouse and sprayed with 2% insecticidal soap to
114 remove pests. Once inside the greenhouse, the plant was left to mature and set seed autonomously via
115 self-pollination. Mature siliques were harvested in July 2013 and seeds were stored in paper envelopes
116 at 4 °C.

117 *S*₂ generation

118 In May 2016, we removed a subset of 10 seeds of the *S*₁ generation from cold storage and germinated in
119 a 60mm × 15mm petri dish containing filter paper covered with a mixture of autoclaved ProMix soil and
120 silica sand (1:9 ratio). We added distilled water until saturation and thereafter petri dishes were sealed
121 with paraffin wax before storing in the dark at 4° C. Of these, six seedlings germinated and were
122 retained for transcriptome sequencing, divided into one of two treatments. The first true leaf from each
123 of the three plants in the experimental treatment were cut with scissors. We used a Kimwipe tissue
124 saturated with either 0.4mM jasmonic acid (JA) dissolved in 10% ethanol (treatment) or 10% ethanol
125 alone (control), adhered directly to maintain contact the cut site (treatment) or uncut leaf (control). We
126 replaced the saturated Kimwipe every 8h to maintain the signal. After 48h of treatment, we harvested
127 seedlings and preserved treated leaves and root tissue in liquid nitrogen, to be used for RNA purification
128 and transcriptome sequencing.

129 DNA Isolation & Library Construction

130 In September 2013, frozen tissue from the *S*₁ genotype was pulverized and extracted using a
131 Cetyltrimethyl Ammonium Bromide (CTAB) protocol (Clarke, 2009) with the following modifications.
132 After pulverising tissue in a bead mill homogenizer at 60 Hz for 60s, we added 1mL chilled wash buffer
133 (200mM Tris-HCl pH 8.0, 50mM EDTA, 250mM NaCl) and incubated on ice for 10min. The purpose of
134 this wash step is to remove secondary metabolites after disruption of cell walls but prior to cell lysis with
135 CTAB. Following this initial wash step, we spun tubes in a microcentrifuge at 4000 g and 4° C for 10
136 minutes, then discarded the supernatant and added another 1mL of wash buffer. This was repeated
137 once more for a total of three wash cycles until no coloration was visible in the supernatant. After final
138 discard of the supernatant and addition of warm lysis buffer as per the CTAB protocol, we vortexed
139 tubes briefly to resuspend plant cells. After completion of the CTAB protocol, pellets were dissolved in
140 50 µL of reverse osmosis (RO) H₂O and sent to Centre d'expertise et de services Génome Québec
141 (Génome Québec) for library preparation and sequencing.

142 We used four separate sequencing libraries for genome assembly: (i) One whole-genome shotgun
143 sequencing library using the Illumina TruSeq DNA v1 preparation kit with a target fragment length of
144 150b. (ii & iii) Two Illumina Nextera MatePair libraries with target lengths of 5Kb and 10Kb. These three
145 libraries were multiplexed and sequenced on a single flowcell of Illumina HiSeq 2000 using 2× 100b
146 paired end (PE) sequencing chemistry. (iv) Target fragment lengths of 450b using the Illumina TruSeq
147 DNA v1 and sequenced on Illumina MiSeq with 2× 250b paired-end reads.

148 RNA Isolation & Library Construction

149 For RNA purification, we pulverized frozen leaf and root tissue in March 2017, in the same manner as
150 the DNA extraction protocol outlined in the previous section. After pulverizing the tissue, we extracted
151 whole RNA from each plant separately using Invitrogen's TRIzol reagent, following the manufacturer's
152 protocol (Pub No. MAN0001271 Rev. A.0). We sequenced four of the six extractions with the highest
153 RNA yields at Génome Québec using the Illumina TruSeq LT kit and multiplex kit for sequencing on a
154 single lane of Illumina MiSeq with 2× 125b paired-end reads. Four separate libraries were sequenced
155 based on tissue and treatment: Control Leaf (CL), Control Root (CR), Treated Leaf (TL), and Treated
156 Root (TR).

157 Data Processing Methods

158 Raw sequencing data was processed and demultiplexed by Génome Québec, and copied to the
159 Frontenac cluster hosted by the Centre for Advanced Computing (CAC) at Queen's University and the
160 Cedar cluster maintained by Simon Fraser University on behalf of Compute Canada. The CAC maintains
161 the Rosalind Franklin Cluster for Analysis of Complex Genomes, which is a 256-core computing cluster
162 with 2 TB of RAM. We used this hardware for the memory-intensive steps *de novo* genome assembly,
163 with the remaining analyses completed using shared Frontenac and Cedar clusters, and on personal
164 computers. All FASTQ files from both experiments passed quality controls using **fastqc** (version 0.11.5)
165 (Andrews, 2010). We used the raw, demultiplexed FASTQ files for *de novo* assembly, but the
166 transcriptome data were pre-processing using **cutadapt** (Martin, 2011, p. 201) to trim adapters and
167 removing quality reads shorter than 25b prior to assembly.

168 Our genome assembly pipeline involved two main steps. First, we used **ALLPATHS-LG** (Gnerre et al.,
169 2011) version R52488 to assemble contigs from both the HiSeq and MiSeq paired-end libraries and then
170 to link contigs into scaffolds using the 5Kb and 10Kb MatePair libraries. The analysis parameters
171 included PLOIDY = 2 and HAPLOIDIFY = TRUE to perform a diploid genome assembly. Although our
172 genome is likely hexaploidy, polyploid models are not supported by **ALLPATHS-LG** and low
173 heterozygosity is expected given the high selfing rates in natural populations. Second, we joined
174 scaffolds from **ALLPATHS-LG** into larger mega-scaffolds using **redundans** (Pryszcz and Gabaldón, 2016)
175 version 0.13c, with the following parameters: *identity* = 0.9, *iters* = 5, *joins* = 5, *limit* = 1, *linkratio* = 0.7,
176 *mapq* = 10, *minlength* = 1000 and *overlap* = 0.75. We repeated this script four times with output
177 scaffolds of the prior run acting as input scaffolds given the long run-time required (~28d). This second
178 combines scaffolds with overlapping similarity, resulting in mega-scaffolds that can span across multiple
179 chromosomes.

180 Sequencing data from the transcriptome experiment were cleaned and assembled with **trinity** (Grabherr
181 et al., 2011) following protocols outlined on the software documentation and in Haas *et al.* (Haas et al.,
182 2013). We used default parameters and the quality of the assembly was analyzed using the custom perl
183 scripts included in the **trinity** package to examine full length transcripts and Contig Nx lengths (i.e.

184 *analyze_blastPlus_topHit_coverage.pl*, *TrinityStats.pl*). Additionally, we mapped read pairs to the
185 transcriptome assembly to assess read content using **bowtie2** (version 2.3.3.1) with default parameters
186 (Langmead and Salzberg, 2012). We used **Transdecoder** to predict open reading frames in transcripts
187 before using **Trinotate** to annotate and analyze assembled transcripts (Haas et al., 2013).

188 As a first step in annotation, we established a detailed repeat library. Miniature Inverted Transposable
189 Elements (MITES) represent the most abundant transposable elements in plants and were identified
190 using **MITE Tracker** (Crescente et al., 2018). Long Terminal Repeat (LTR) elements are less common but
191 occupy a larger proportion of the genome, and were identified using the **GenomeTools** package
192 (Gremme et al., 2013). To reduce the number of false positive LTR transposons, only those that
193 contained PPT (poly purine extract) or PBS (primer binding sites) were kept and the rest filtered. We
194 further filtered the LTR candidates to eliminate three main sources of false positives: tandem local
195 repeats such as centromeric repeats, local gene clusters derived from recent gene duplications, and two
196 other transposable elements located in adjacent regions. We also identified elements with nested
197 insertions. After processing known MITES and LTR elements, we identified additional repetitive
198 sequences using **RepeatModeler** against a transposase database and excluded gene fragments using
199 **ProtExcluder**.

200 The annotation pipeline **maker** (Campbell et al. 2014) was used to identify gene models and predict
201 functional annotations in the draft genome. Both *est2genome* and *protein2genome* modes were used
202 initially to make *ab initio* gene predictions from EST and protein evidence, respectively. The EST
203 evidence was based on our own transcriptome data whereas the protein evidence was gathered from
204 the reference proteomes of six closely related plants available from the SwissProt database: *Arabidopsis*
205 *thaliana*, *Glycine max*, *Brassica oleracea*, *Medicago truncatula*, *Brassica napus*, and *Brassica rapa*. From
206 the first round of annotation, high confidence models were predicted by the *maker2zff* command with
207 default minimums (50% of splice sites confirmed by EST alignment, 50% of exons match an EST
208 alignment, 50% of exons overlap any evidence, and maximum AED of 0.5). These predictions were used
209 to train the *ab initio* gene predictor **snap** (Korf, 2004). A second round of **maker** was run using the
210 hidden Markov model (HMM) from **snap** rather than the *est2genome* mode. All other settings were the
211 same as for the first run, with the transcripts now being used only as evidence to support *ab initio* gene
212 predictions. Two more rounds of annotation and gene prediction improvement followed. An **Augustus**
213 gene prediction file was also generated for use as a second *ab initio* prediction (Stanke and
214 Morgenstern, 2005). For the final round of annotation in **maker**, the final HMM file from **snap**, the *A.*
215 *petiolata* species **Augustus** library and gene prediction file were all used in addition to the following
216 settings: always complete, single exons, and using correct EST fusion.

217 We used Benchmarking Universal Single-Copy Orthologs (**BUSCO**) v4.1 (Simão et al., 2015) to assess the
218 assembly quality of the final draft genome. We used plant lineage delineation from the EmbryophytaDB
219 V10 database, focusing on universal orthologs present in >90% of lineages, resulting in a total of 1,440
220 BUSCO orthologs. We used **minimap2** check for assembly contiguity and synteny with the model plant
221 *Arabidopsis thaliana* using the TAIR 10 assembly with up to 20% sequence divergence. Plots generated
222 with the R package **pafr**.

223 The genome annotation files were then curated through **deFusion** to resolve fused gene annotation
224 problems (Wang, 2020), as well as **EvidenceModeler** (Haas et al., 2008) to combine *ab initio* gene
225 predictions and protein and transcript alignments into weighted consensus gene structures. The

226 functional annotations were then created using NCBI BLAST+ and InterProScan (Jones et al., 2014) by
227 adding new names, domains, and putative functions to improve the utility of the genome database.

228 Data Availability

229 Raw data used for genome assembly, transcriptome assembly, and the final draft genome can be found
230 in the NCBI SRA database (accession numbers TBD) under BioProject **SUB9096116** (BioSample Accession:
231 **SAMN17958863**). Analysis scripts are available on GitHub
232 (https://github.com/turkrhen/snapping_turtle_genome_scripts).

233 Results and Discussion

234 Whole Genome Shotgun (WGS) sequencing of genomic DNA produced 45.8Gb from 229 million HiSeq
235 reads and an additional 7.9Gb from 15.8 million MiSeq reads. This represents an estimated 68× average
236 coverage of the genome with an estimated genome size of 1.07Gb. This is similar to published size of
237 1.35Gb estimated by flow cytometry (Barow and Meister, 2002). Initial assembly with **ALLPATHS-LG**
238 resulted in 16,743 contigs longer than 1Kb. Technically, these are scaffolded contigs linked using mate
239 pairs, however we refer to these as contigs to avoid confusion with the scaffolds created by merging
240 heterozygous loci in the **redundans** program. The final assembled genome was 1.08 Gb long across 694
241 scaffolds larger than 1Kb, and more than 50% of the genome is contained in the ten longest scaffolds
242 (Table 1).

243 Cytological studies of *A. petiolata* suggest variation in chromosome number and ploidy (Cavers et al.,
244 1979; Weiss-Schneeweiss and Schneeweiss, 2003), and whether the hexaploids are allo- or auto-
245 polyploids. Approximately 36% of contigs (6,045 of 16,743) remain heterozygous after **redundans**
246 assembly; however, average identity within these contigs was 94.8%, meaning that an average of just
247 2% of sites in the genome assembly are heterozygous. Similarly, just 3.2% of **BUSCO** genes mapped to
248 more than two scaffolds (46 of 1440), consistent with a low level of heterozygosity and minimal
249 duplication of housekeeping genes.

250 This relatively low level of heterozygosity is consistent with a diploid or highly-inbred autopolyploid. A
251 relatively simple genome combined with geographic variation in ploidy make *A. petiolata* an appealing
252 species to study the role of polyploidy in local adaptation and range expansion, which is an active area
253 of research (e.g. (Pandit et al., 2011; Payseur and Rieseberg, 2016; te Beest et al., 2012)).

254 **Table 1.** Assembly statistics for the *Alliaria petiolata* genome

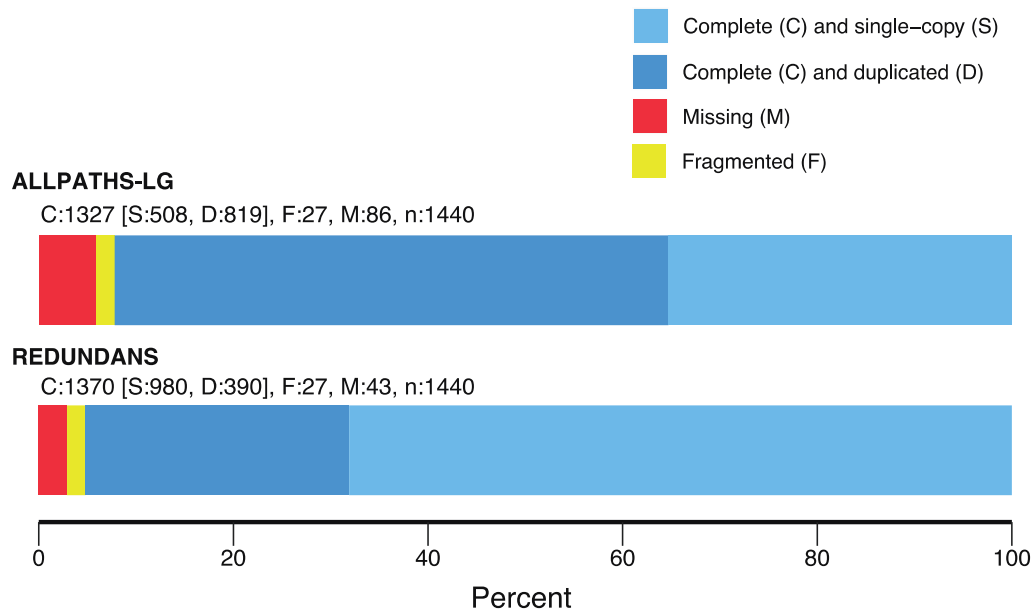
Statistic	Value
# contigs (>= 1000b)	694
# contigs (>= 50,000b)	227
Total length	1,075,010,735
Total length (>= 50,000b)	1,071,536,925
Largest contig	485,611,451
GC (%)	37.2
N ₅₀	121,941,980
N ₇₅	40,840,077
L ₅₀	2

L₇₅ 5
Mean Sequence Length 1,549,006.82

255

256 Most (98.7%) of the essential single-copy genes from **BUSCO** mapped to our assembly (Fig. 1), with
257 71.5% occurring only once in the assembled genome. Similarly, 98% of sequence reads from the
258 transcriptome experiment mapped successfully to the assembled genome. Sequencing of the
259 transcriptome libraries yielded a total of 68.1 Gb from 272.5 million paired reads. Trimming sequence
260 reads for quality reduced usable data by less than 2%.

261



262

263 **Figure 1.** Percentage of predicted single-copy plant genes from **BUSCO** that are
264 found one (light blue) or more times (dark blue), or are missing (red) or fragmented
265 (yellow) in the annotated genome assembly of *Alliaria petiolata*.

266 Our *de novo* transcriptome assembly included 699,048 putative isoform “transcripts”
267 representing 535 Mb with N50 of 1,233 base pairs. The minimum transcript length was 201 as set
268 by the Trinity default parameter while 1382 (~1.98%) of transcripts were longer than 5Kb. These
269 transcripts clustered into 350,672 hypothetical genes with an average of 2.18 isoforms and 26,910
270 (~7.67%) of hypothetical genes having more than five isoforms. A BLAST search of hypothetical genes to
271 the SwissProt protein database matched 10,352 proteins with at least 90% coverage of the query
272 sequence, including 7,930 proteins with 100% coverage (Table 2).

273 Our TE annotation analysis identified 8,220 unique sequences across the *A. petiolata* genome. Of these,
274 112 were classified as LTRs with relatively recent origin (99% similarity), 240 as relatively old (85%), and
275 7,137 were classified as miniature inverted transposable elements (MITE). An additional 731 sequences
276 were found to match the DNA transposase family. After masking TE sequences the final gene set from
277 **maker** included 64,770 gene predictions with an average of 6 exons and an average exon length of 251b
278 (Table 2).

279 A dot-plot comparison of gene synteny with the model plant *Arabidopsis thaliana* (TAIR 10) revealed
280 large blocks of orthologous sequence (Figure 2). However, the arrangement of synteny blocks shows a
281 complete re-arrangement of *A. thaliana* chromosomes when mapped to the *A. petiolata* contigs.
282 Despite a high level of asynteny that is characteristic of the Brassicaceae family, the conservation of
283 large synteny blocks can help to identify candidate genes and genetic loci of interest for understanding
284 plant invasions. Future research could also investigate whether gene rearrangements occur among
285 geographically and historically isolated populations of *A. petiolata*, and whether this genomic
286 architecture has played an important role in the spread of the species.

287

288 **Table 2.** Summary statistics of genes annotated for the *Alliaria petiolata* genome assembly.

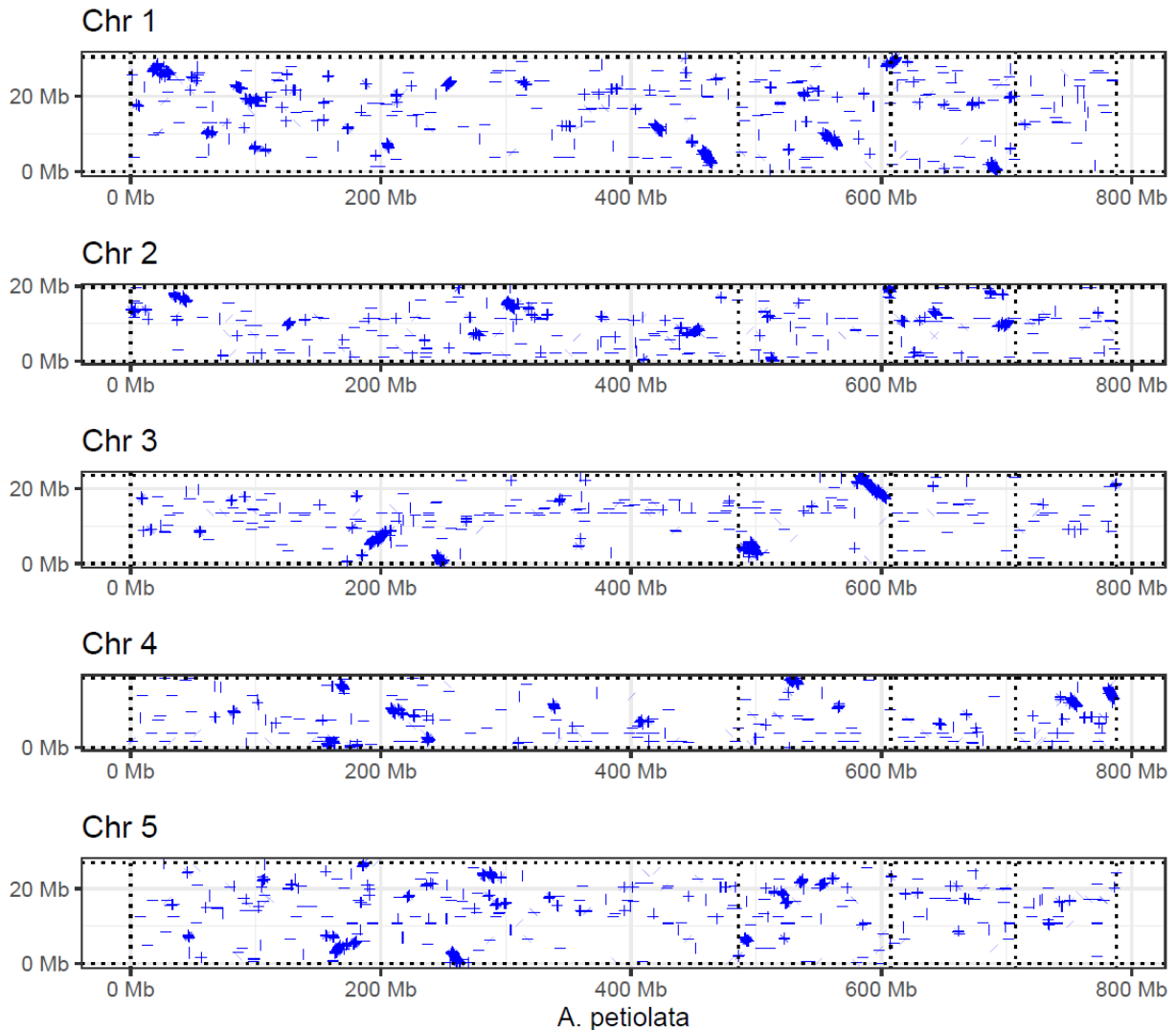
Statistic	Value
Number of genes	64,770
Number of exons	408,155
Number of introns	343,385
Overlapping genes	9,669
Contained genes	1,464
Total gene length	210,804,785
Total exon length	102,316,121
Total intron length	109,175,434
Total CDS length	842,788
% of genome covered by genes	19.6
% of genome covered by CDS	7.8
Mean mRNAs per gene	1
Mean exons per mRNA	6
Mean introns per mRNA	5
Mean gene length	3,255
Mean intron length	318
Mean exon length	251
Mean CDS length	1301

289

290

291

292



293

294 **Figure 2.** Dot-plot showing blocks of synteny between scaffolds of the *Alliaria petiolata* assembly (x-
295 axis), and five chromosomes of the model plant *Arabidopsis thaliana*. Blue lines show aligned sequences
296 with up to 20% divergence. Vertical dotted lines denote separation of the major scaffolds of the *A.*
297 *petiolata* assembly.

298 **Conclusions**

299 There is a growing interest in the genetic causes and consequences of range expansion and biological
300 invasion. The field of invasion genetics has emerged from ecological and evolutionary studies of invasive
301 species but lacks well-developed model systems. The draft genome and gene annotation reported here
302 represents an important link from the many field and experimental studies of *A. petiolata* to the genetic
303 architecture of adaptation and invasion. High levels of self-fertility and the resultant low levels
304 heterozygosity observed in the genome will be beneficial for future projects linking ecologically
305 important phenotypes to specific genes. The genomic resources reported here complement available
306 seed resources, experimental findings, and field data to accelerate genomic and molecular studies of *A.*
307 *petiolata* as a candidate for a model system in invasion genetics.

308 Acknowledgements

309 The authors are grateful for bioinformatics support from C Grassa, J Stafford, and H Schmider, hardware
310 support from C MacPhee, and wetlab assistance from M Todesco, W Chen and A Siew. We also thank M
311 Todesco and D Galanti for comments that improved the manuscript. This work was supported by NSERC
312 Discovery grants to RIC and LHR.

313 Literature Cited

- 314 Andrews, S., 2010. FastQC A Quality Control tool for High Throughput Sequence Data [Online] [WWW
315 Document]. URL <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/> (accessed
316 11.19.20).
- 317 Anthony, M.A., Frey, S.D., Stinson, K.A., 2017. Fungal community homogenization, shift in dominant
318 trophic guild, and appearance of novel taxa with biotic invasion. *Ecosphere* 8, e01951.
319 <https://doi.org/10.1002/ecs2.1951>
- 320 Barow, M., Meister, A., 2002. Lack of correlation between AT frequency and genome size in higher
321 plants and the effect of nonrandomness of base sequences on dye binding. *Cytometry* 47, 1–7.
322 <https://doi.org/10.1002/cyto.10030>
- 323 Barrett, S.C.H., 2015. Foundations of invasion genetics: The Baker and Stebbins legacy. *Molecular*
324 *Ecology* 24, 1927–1941. <https://doi.org/10.1111/mec.13014>
- 325 Barrett, S.C.H., Colautti, R.I., Dlugosch, K.M., Rieseberg, L.H., 2016. *Invasion Genetics: The Baker and*
326 *Stebbins Legacy*. Wiley.
- 327 Bock, D.G., Caseys, C., Cousens, R.D., Hahn, M.A., Heredia, S.M., Huebner, S., Turner, K.G., Whitney, K.D.,
328 Rieseberg, L.H., 2015. What we still don't know about invasion genetics. *Molecular Ecology* 24,
329 2277–2297. <https://doi.org/10.1111/mec.13032>
- 330 Boheemen, L.A. van, Lombaert, E., Nurkowski, K.A., Gauffre, B., Rieseberg, L.H., Hodgins, K.A., 2017.
331 Multiple introductions, admixture and bridgehead invasion characterize the introduction history
332 of *Ambrosia artemisiifolia* in Europe and Australia. *Molecular Ecology* 26, 5421–5434.
333 <https://doi.org/10.1111/mec.14293>
- 334 Bossdorf, O., Prati, D., Auge, H., Schmid, B., 2004. Reduced competitive ability in an invasive plant.
335 *Ecology Letters* 7, 346–353. <https://doi.org/10.1111/j.1461-0248.2004.00583.x>
- 336 Bourne, S.D., Hudson, J., Holman, L.E., Rius, M., 2020. Marine Invasion Genomics: Revealing Ecological
337 and Evolutionary Consequences of Biological Invasions, in: Oleksiak, M.F., Rajora, O.P. (Eds.),
338 *Population Genomics: Marine Organisms, Population Genomics*. Springer International
339 Publishing, Cham, pp. 363–398. https://doi.org/10.1007/13836_2018_21
- 340 Callaway, R.M., Cipollini, D., Barto, K., Thelen, G.C., Hallett, S.G., Prati, D., Stinson, K., Klironomos, J.,
341 2008. Novel weapons: Invasive plant suppresses fungal mutualists in America but not in its
342 native Europe. *Ecology* 89, 1043–1055. <https://doi.org/10.1890/07-0370.1>
- 343 Cavers, P.B., Heagy, M.I., Kokron, R.F., 1979. The Biology of Canadian Weeds: 35. *Alliaria petiolata* (M.
344 Bieb.) Cavara and Grande. *Canadian Journal of Plant Science* 59, 217–229.
- 345 Cipollini, D., Cipollini, K., 2016. A review of garlic mustard (*Alliaria petiolata*, Brassicaceae) as an
346 allelopathic plant. *tbot* 143, 339–348. <https://doi.org/10.3159/TORREY-D-15-00059>
- 347 Clarke, J.D., 2009. Cetyltrimethyl Ammonium Bromide (CTAB) DNA Miniprep for Plant DNA Isolation.
348 Cold Spring Harb Protoc 2009, pdb.prot5177. <https://doi.org/10.1101/pdb.prot5177>
- 349 Colautti, R., Franks, S.J., Hufbauer, R.A., Kotanen, P.M., Torchin, M., Byers, J.E., Pyšek, P., Bossdorf, O.,
350 2014. The Global Garlic Mustard Field Survey (GGMFS): challenges and opportunities of a
351 unique, large-scale collaboration for invasion biology. *NeoBiota* 21, 29–47.
352 <https://doi.org/10.3897/neobiota.21.5242>

- 353 Crescente, J.M., Zavallo, D., Helguera, M., Vanzetti, L.S., 2018. MITE Tracker: an accurate approach to
354 identify miniature inverted-repeat transposable elements in large genomes. *BMC Bioinformatics*
355 19, 348. <https://doi.org/10.1186/s12859-018-2376-y>
- 356 Dlugosch, K.M., Parker, I.M., 2008. Invading populations of an ornamental shrub show rapid life history
357 evolution despite genetic bottlenecks. *Ecol. Lett.* 11, 701–709. <https://doi.org/10.1111/j.1461-0248.2008.01181.x>
- 359 Duchesneau, K., Golemic, A., Colautti, R.I., Antunes, P.M., 2020. Functional shifts of soil microbial
360 communities associated with *Alliaria petiolata* invasion. *bioRxiv* 2020.08.13.248849.
361 <https://doi.org/10.1101/2020.08.13.248849>
- 362 Esmailbegi, S., Al-Shehbaz, I.A., Pouch, M., Mandáková, T., Mummenhoff, K., Rahiminejad, M.R.,
363 Mirtadzadini, M., Lysak, M.A., 2018. Phylogeny and systematics of the tribe Thlaspidae
364 (Brassicaceae) and the recognition of two new genera. *TAXON* 67, 324–340.
365 <https://doi.org/10.12705/672.4>
- 366 Gnerre, S., MacCallum, I., Przybylski, D., Ribeiro, F.J., Burton, J.N., Walker, B.J., Sharpe, T., Hall, G., Shea,
367 T.P., Sykes, S., Berlin, A.M., Aird, D., Costello, M., Daza, R., Williams, L., Nicol, R., Gnirke, A.,
368 Nusbaum, C., Lander, E.S., Jaffe, D.B., 2011. High-quality draft assemblies of mammalian
369 genomes from massively parallel sequence data. *PNAS* 108, 1513–1518.
370 <https://doi.org/10.1073/pnas.1017351108>
- 371 Grabherr, M.G., Haas, B.J., Yassour, M., Levin, J.Z., Thompson, D.A., Amit, I., Adiconis, X., Fan, L.,
372 Raychowdhury, R., Zeng, Q., Chen, Z., Mauceli, E., Hacohen, N., Gnirke, A., Rhind, N., di Palma,
373 F., Birren, B.W., Nusbaum, C., Lindblad-Toh, K., Friedman, N., Regev, A., 2011. Trinity:
374 reconstructing a full-length transcriptome without a genome from RNA-Seq data. *Nat Biotechnol*
375 29, 644–652. <https://doi.org/10.1038/nbt.1883>
- 376 Gremme, G., Steinbiss, S., Kurtz, S., 2013. GenomeTools: A Comprehensive Software Library for Efficient
377 Processing of Structured Genome Annotations. *IEEE/ACM Transactions on Computational*
378 *Biology and Bioinformatics* 10, 645–656. <https://doi.org/10.1109/TCBB.2013.68>
- 379 Haas, B.J., Papanicolaou, A., Yassour, M., Grabherr, M., Blood, P.D., Bowden, J., Couger, M.B., Eccles, D.,
380 Li, B., Lieber, M., MacManes, M.D., Ott, M., Orvis, J., Pochet, N., Strozzi, F., Weeks, N.,
381 Westerman, R., William, T., Dewey, C.N., Henschel, R., LeDuc, R.D., Friedman, N., Regev, A.,
382 2013. De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for
383 reference generation and analysis. *Nature Protocols* 8, 1494–1512.
384 <https://doi.org/10.1038/nprot.2013.084>
- 385 Haas, B.J., Salzberg, S.L., Zhu, W., Pertea, M., Allen, J.E., Orvis, J., White, O., Buell, C.R., Wortman, J.R.,
386 2008. Automated eukaryotic gene structure annotation using EvidenceModeler and the
387 Program to Assemble Spliced Alignments. *Genome Biology* 9, R7. <https://doi.org/10.1186/gb-2008-9-1-r7>
- 389 Haribal, M., Yang, Z., Attygalle, A.B., Renwick, J.A.A., Meinwald, J., 2001. A Cyanoallyl Glucoside from
390 *Alliaria petiolata*, as a Feeding Deterrent for Larvae of *Pieris napi* oleracea. *J. Nat. Prod.* 64, 440–
391 443. <https://doi.org/10.1021/np000534d>
- 392 Jones, P., Binns, D., Chang, H.-Y., Fraser, M., Li, W., McAnulla, C., McWilliam, H., Maslen, J., Mitchell, A.,
393 Nuka, G., Pesseat, S., Quinn, A.F., Sangrador-Vegas, A., Scheremetjew, M., Yong, S.-Y., Lopez, R.,
394 Hunter, S., 2014. InterProScan 5: genome-scale protein function classification. *Bioinformatics* 30,
395 1236–1240. <https://doi.org/10.1093/bioinformatics/btu031>
- 396 Korf, I., 2004. Gene finding in novel genomes. *BMC Bioinformatics* 5, 59. <https://doi.org/10.1186/1471-2105-5-59>
- 398 Langmead, B., Salzberg, S.L., 2012. Fast gapped-read alignment with Bowtie 2. *Nature Methods* 9, 357–
399 359. <https://doi.org/10.1038/nmeth.1923>

- 400 Lankau, R.A., Nuzzo, V., Spyreas, G., Davis, A.S., 2009. Evolutionary limits ameliorate the negative impact
401 of an invasive plant. *Proc. Natl. Acad. Sci. U. S. A.* 106, 15362–15367.
402 <https://doi.org/10.1073/pnas.0905446106>
- 403 Lewis, K.C., Bazzaz, F.A., Liao, Q., Orians, C.M., 2006. Geographic patterns of herbivory and resource
404 allocation to defense, growth, and reproduction in an invasive biennial, *Alliaria petiolata*.
405 *Oecologia* 148, 384–395. <https://doi.org/10.1007/s00442-006-0380-9>
- 406 Mack, R.N., Simberloff, D., Lonsdale, W.M., Evans, H., Clout, M., Bazzaz, F.A., 2000. Biotic invasions:
407 Causes, epidemiology, global consequences, and control. *Ecological Applications* 10, 689–710.
- 408 Martin, M., 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads.
409 *EMBnet.journal* 17, 10–12. <https://doi.org/10.14806/ej.17.1.200>
- 410 Mooney, H.A., Cleland, E.E., 2001. The evolutionary impact of invasive species. *PNAS* 98, 5446–5451.
411 <https://doi.org/10.1073/pnas.091093398>
- 412 Nuzzo, V., 1993. Distribution and spread of the invasive biennial *Alliaria petiolata* (garlic mustard) in
413 North America., in: McKnight, B. (Ed.), *Biological Pollution: The Control and Impact of Invasive*
414 *Exotic Species*. Indiana Academy of Science, Indianapolis, Indiana, pp. 137–145.
- 415 Pandit, M.K., Pockock, M.J.O., Kunin, W.E., 2011. Ploidy influences rarity and invasiveness in plants.
416 *Journal of Ecology* 99, 1108–1115. <https://doi.org/10.1111/j.1365-2745.2011.01838.x>
- 417 Payseur, B.A., Rieseberg, L.H., 2016. A genomic perspective on hybridization and speciation. *Molecular*
418 *Ecology* 25, 2337–2360. <https://doi.org/10.1111/mec.13557>
- 419 Prati, D., Bossdorf, O., 2004. Allelopathic inhibition of germination by *Alliaria petiolata* (Brassicaceae).
420 *American Journal of Botany* 91, 285–288. <https://doi.org/10.3732/ajb.91.2.285>
- 421 Pryszcz, L.P., Gabaldón, T., 2016. Redundans: an assembly pipeline for highly heterozygous genomes.
422 *Nucleic Acids Res* 44, e113–e113. <https://doi.org/10.1093/nar/gkw294>
- 423 Sax, D.F., Stachowicz, J.J., Gaines, S.D., 2005. *Species Invasions: Insights into Ecology, Evolution, and*
424 *Biogeography*. Sinauer.
- 425 Simão, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V., Zdobnov, E.M., 2015. BUSCO: assessing
426 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 31,
427 3210–3212. <https://doi.org/10.1093/bioinformatics/btv351>
- 428 Stanke, M., Morgenstern, B., 2005. AUGUSTUS: a web server for gene prediction in eukaryotes that
429 allows user-defined constraints. *Nucleic Acids Res* 33, W465–W467.
430 <https://doi.org/10.1093/nar/gki458>
- 431 Stinson, K., Kaufman, S., Durbin, L., Lowenstein, F., 2007. Impacts of garlic mustard invasion on a forest
432 understory community. *Conserv Biol* 14, 73–88. [https://doi.org/10.1656/1092-6194\(2007\)14\[73:IOGMIO\]2.0.CO;2](https://doi.org/10.1656/1092-6194(2007)14[73:IOGMIO]2.0.CO;2)
- 433
- 434 te Beest, M., Le Roux, J.J., Richardson, D.M., Brysting, A.K., Suda, J., Kubešová, M., Pyšek, P., 2012. The
435 more the better? The role of polyploidy in facilitating plant invasions. *Ann Bot* 109, 19–45.
436 <https://doi.org/10.1093/aob/mcr277>
- 437 USDA, N., 2020. The PLANTS Database (<http://plants.usda.gov>, 10 April 2020). National Plant Data Team,
438 Greensboro, NC.
- 439 Wang, J., 2020. *wjidea/defusion*.
- 440 Weiss-Schneeweiss, H., Schneeweiss, G.M., 2003. Karyological investigations of selected angiosperms
441 from Georgia and Azerbaijan. *Acta Biologica Cracoviensia Series Botanica* 45, 49–56.
442