

# Genome Improvement and Genetic Map Construction for *Aethionema arabicum*, the First Divergent Branch in the Brassicaceae Family

Thu-Phuong Nguyen,<sup>\*</sup> Cornelia Mühlich,<sup>†</sup> Setareh Mohammadin,<sup>\*</sup> Erik van den Bergh,<sup>\*</sup>

Adrian E. Platts,<sup>‡</sup> Fabian B. Haas,<sup>†</sup> Stefan A. Rensing,<sup>\*,§,1</sup> and M. Eric Schranz<sup>\*,1</sup>

<sup>\*</sup>Biosystematics Group, Wageningen University, Droevendaalsesteeg 1, 6708 PB Wageningen, The Netherlands, <sup>†</sup>Faculty of Biology, Plant Cell Biology, University of Marburg, Karl-von-Frisch-Str. 8, 35043 Marburg, Germany, <sup>‡</sup>Center for Genomics and Systems Biology, Department of Biology, New York University, New York, and <sup>§</sup>LOEWE Center for Synthetic Microbiology (SYNMIKRO), Philipps University of Marburg, Germany

ORCID IDs: 0000-0001-9865-574X (E.v.d.B.); 0000-0002-7711-5282 (F.B.H.); 0000-0002-0225-873X (S.A.R.); 0000-0001-6777-6565 (M.E.S.)

**ABSTRACT** The genus *Aethionema* is a sister-group to the core-group of the Brassicaceae family that includes *Arabidopsis thaliana* and the Brassica crops. Thus, *Aethionema* is phylogenetically well-placed for the investigation and understanding of genome and trait evolution across the family. We aimed to improve the quality of the reference genome draft version of the annual species *Aethionema arabicum*. Second, we constructed the first *Ae. arabicum* genetic map. The improved reference genome and genetic map enabled the development of each other. We started with the initially published genome (version 2.5). PacBio and MinION sequencing together with genetic map v2.5 were incorporated to produce the new reference genome v3.0. The improved genome contains 203 MB of sequence, with approximately 94% of the assembly made up of called (non-gap) bases, assembled into 2,883 scaffolds (with only 6% of the genome made up of non-called bases (Ns)). The N<sub>50</sub> (10.3 MB) represents an 80-fold increase over the initial genome release. We generated a Recombinant Inbred Line (RIL) population that was derived from two ecotypes: Cyprus and Turkey (the reference genotype). Using a Genotyping by Sequencing (GBS) approach, we generated a high-density genetic map with 749 (v2.5) and then 632 SNPs (v3.0) was generated. The genetic map and reference genome were integrated, thus greatly improving the scaffolding of the reference genome into 11 linkage groups. We show that long-read sequencing data and genetics are complementary, resulting in an improved genome assembly in *Ae. arabicum*. They will facilitate comparative genetic mapping work for the Brassicaceae family and are also valuable resources to investigate wide range of life history traits in *Aethionema*.

## KEYWORDS

*Aethionema arabicum*  
Brassicaceae  
genome  
improvement  
genetic map  
PacBio  
MinION  
Genotyping by Sequencing

## BACKGROUND

The genus *Aethionema* belongs to the important plant family Brassicaceae. The crucifers contain many species of interest, such

as the Brassica crop plants (e.g., *B. rapa*, *B. oleracea* and *B. napus*), ornamental plants (such as the genera *Aubrieta*, *Iberis*, *Lunaria* and *Draba*) and research model plant species (including *Arabidopsis thaliana*, *A. lyrata*, *Capsella rubella* and *Arabis alpina*). Phylogenetic studies have established *Aethionema* as the sister-group of the core-group in the family (Beilstein *et al.* 2008; Huang *et al.* 2016; Guo *et al.* 2017). Thus, *Aethionema* holds an essential phylogenetic position for studies on genome and trait evolution across the Brassicaceae family.

The monogeneric tribe *Aethionemeae* consists of 57 species and is distributed mainly in the Irano-Turanian region, a hot spot for species radiation and speciation (Al-Shehbaz *et al.* 2006; Franzke *et al.* 2011). This tribe displays various interesting morphological and ecological characteristics, especially fruit and

Copyright © 2019 Nguyen *et al.*

doi: <https://doi.org/10.1534/g3.119.400657>

Manuscript received March 29, 2019; accepted for publication September 11, 2019; published Early Online September 25, 2019.

This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Supplemental material available at FigShare: <https://doi.org/10.25387/g3.8233055>.

<sup>1</sup>Corresponding authors: E-mail: stefan.rensing@biologie.uni-marburg.de, eric.schranz@wur.nl

seed heteromorphism. Heteromorphism is defined as the production of two or more distinct fruit or seed morphs on the same individual (Imbert 2002), which includes morphological size, shape and color; physiological dormancy and germination of fruits and seeds. *Aethionema arabicum* is one of the seven reported heteromorphic species of *Aethionema* (Lenser *et al.* 2016; Mohammadin *et al.* 2017). *Aethionema arabicum* is a small diploid annual, with a short life cycle starting from seed germination to the end of the vegetative development in spring, followed by reproduction and the end of life cycle in summer (Bibalani 2012). Both annual life history and heteromorphism probably evolved as adaptive responses to unpredictable environments, especially dry arid habitats, indicating a wide range of natural variation for ecologically adaptive traits in *Ae. arabicum*.

Owing to its unique phylogenetic position and interesting characteristics, *Ae. arabicum* is an ideal sister-group model for research. Therefore, *Aethionema* genome and genetic resources are desirable. The initially published *Ae. arabicum* draft genome (v1.0) contains 59,101 scaffolds with an N50 of 115,195 bp while the genome was predicted to be 320 Mbp in size with  $n = 11$  (Haudry *et al.* 2013). Here we first aimed to (i) improve the quality of the reference genome and (ii) to construct the first *Ae. arabicum* genetic map. A higher quality version of the genome assembled by the VEGI consortium was later released as version 2.5, which is used as the starting point of our analyses.

High throughput sequencing using Pacific Biosciences (PacBio) and Oxford Nanopore MinION (MinION) technology followed and resolved many uncalled gaps in the v2.5 genome and supported further super-scaffolding, which resulted in genome v3.0.

The genetic map was constructed using Genotyping by Sequencing (GBS) on a Recombinant Inbred Line (RIL) population. The 216 RILs were derived from a cross between Turkey (reference ecotype) and Cyprus ecotypes. The first version of genetic map v2.5 was obtained based on genome v2.5 with 746 Single Nucleotide Polymorphism (SNP) markers. The later genetic map v3.0 was built with 626 SNPs generated based on genome v3.0.

Here we show that the long-read genome assembly and the genetic map of *Ae. arabicum* supported the development of each other. They will serve as a substantial resource for further research on *Aethionema* as well as the Brassicaceae family.

## DATA DESCRIPTION AND METHODS

### Overview of the workflow

An overview of the improvements of the genome of *Aethionema arabicum* and the generation and improvement of its genetic map are depicted in Figure 1. The genome draft version 1.0 was first improved by Ray (Boisvert *et al.* 2010) and AllPathsLGs (Gnerre *et al.* 2011) and led to the release of genome v2.5 (available on genomevolution.org). Genome v2.5 was used as a basis for SNP calling after GBS of the RILs. This generated SNP markers used to construct the genetic map v2.5. Scaffolds were ordered with AllMaps (Tang *et al.* 2015) based on the maximum co-linearity to genome v2.5 and genetic map v2.5. This resulted in genome vAM. Gap filling and super-scaffolding improvement for genome vAM was obtained by PacBio sequencing leading to genome v2.6. PBjelly2 (English *et al.* 2012) run using the MinION reads further improved genome v2.6 to v3.0. We revisited the genetic map v2.5 by recalling SNPs according to genome v3.0 and constructed a genetic map v3.0 with the newly called SNP markers. Below we describe the workflow in detail in the three following sections: (i) the initial

genome assembly, (ii) genetic map construction and (iii) genome improvement.

**The initial genome (v2.5): The starting point:** *Genome re-assembly using AllPathsLG* The version 1.0 assembly generated by the Ray assembler (Boisvert *et al.* 2010) was fragmented *in silico* into a set of artificial overlapping reads, combined with additional paired end and mate pair data (described in (Haudry *et al.* 2013)), and re-assembled using the AllPathsLG assembler (Gnerre *et al.* 2011). This iterative assembly process leveraged additional short read data and AllPathsLG's extensive error correction to minimize the potential for errors in the first assembly to contribute to the second round of assembly. Gap closing was then performed using GapCloser, part of the SOAPdenovo2 package (Luo *et al.* 2012). Gene annotations were lifted over from assembly version 1.0 to version 2.5 using the LiftOver tool from the UCSC Genome Browser tools package (Kent *et al.* 2002).

Genome version 2.5 contains 3,166 scaffolds, has an N50 of 564,741 bp and was published as version 2.5 on <https://genomevolution.org/coge/>.

**Genetic map construction:** *Plant material* Two *Aethionema arabicum* ecotypes were used, Turkey (TUR) and Cyprus (CYP). The TUR accession comes from the living plant collections at the Botanical Garden in Jena, Germany (Botanischer Garten Jena). The seeds for this genotype were derived from a plant in the Botanical Garden in Nancy, France. The CYP ecotype was collected in 2010 near Kato Moni (coordinates UTM WGS 84: 508374 - 3879403) at an altitude 410 m on pillow lava by Charalambos S Christodoulou (Mohammadin *et al.* 2018).

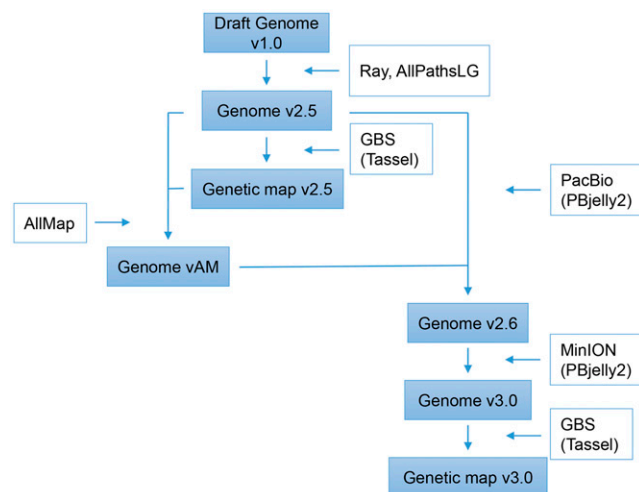
These two ecotypes were used as parents for the development of the recombinant inbred line population, where TUR was the father and CYP the mother. Seeds from initial F<sub>1</sub> plants were used to generate an initial F<sub>2</sub> population. For each of the 216 segregating F<sub>2</sub> plants, a single seed was randomly chosen to further grow and reproduce the next generation. The procedure was repeated until F<sub>8</sub>, when the experiment was performed with 216 RILs.

To grow the plants for the GBS, F<sub>8</sub> seeds of 216 RILs were placed on filter paper, wetted with distilled water, in petri dishes. Imbibed seeds were incubated at 4° in dark for 3 days, followed by germination in the light at 20° for 2 days. Seedlings were transferred to soil pots (10.5 cm diameter 10 cm height) in November 2014. Plants were grown in greenhouse (Wageningen University and Research, the Netherlands) in partially controlled conditions, long day (16 h light and 8 h dark) and at 20°.

### GENOTYPING BY SEQUENCING (GBS)

#### DNA isolation

Young tissues from leaves and flower buds were collected from each F<sub>9</sub> plant for DNA isolation. The DNA isolation was done according to a modified CTAB protocol (Doyle 1991). In brief, plant material was frozen with liquid nitrogen and ground into powder. Each sample was incubated with 500 µl of CTAB buffer at 65° in the water bath for 30 min. After 30 min cooling at room temperature, equal volume (500 µl) of chloroform:isoamylalcohol (24:1 v/v) was added, and vigorously hand-mixed for a min. 400 µl of supernatant was recovered after centrifuging at maximum speed for 5 min. The supernatant was cleaned again with a chloroform:isoamylalcohol step. DNA precipitation was performed by adding an equal volume of cold isopropanol with 30 min incubation on ice and centrifugation at maximum speed for 15 min. The DNA pellet was cleaned twice with 1 ml of 70%



**Figure 1** Overview of the analyses performed in this study. In filled boxes are data sets, approaches and accompanying tools are in open boxes.

ethanol and centrifugation at maximum speed for 5 min. Dry DNA pellet was dissolved in Milli-Q water.

### Constructing GBS libraries

DNA was treated with RNase overnight at 37° with RNase one by Promega. Quality was checked on a 1% agarose gel and DNA quantity was checked with Pico Green. Based on this, DNA was diluted down to 20 ng/μl with MQ water and used in further analysis. GBS was performed in general by following the procedure described in (Elshire *et al.* 2011). Oligonucleotides for creation of common as well as 96 barcoded ApeKI adapters were obtained from Integrated DNA Technologies and diluted to 200 μM. For each barcoded and common adapter, top and bottom strand oligos were combined to a 50 μM annealing molarity in TE to 100 μl total volume. Adapter annealing was carried out in a thermocycler (Applied Biosystems) at 95° for 2 min, ramp to 25° by 0.1 degree per second, hold at 25° for 30 min and 4° forever. Annealed adapters were further diluted to a 0.6 ng/μl concentrated working stock of combined barcoded and common adapter in 96 well microtiter plate and dried using a vacuum oven. For each genomic DNA sample 100 ng (10 ng/μl) was used and added to lyophilized adapter mix and dried down again using a vacuum oven.

Adapter DNA mixtures were digested using 2.5 Units ApeKI (New England Biolabs) for 2 hr at 75° in a 20 μl volume. Digested DNA and Adapters were used in subsequent ligation by 1.6 μl (400 Units/μl) T4DNA Ligase in a 50 μl reaction volume at 22° for one hour followed by heat inactivation at 65° for 30 min. Sets of 96 digested DNA samples, each with a different barcode adapter, were combined (10μl each) and purified using a Qiaquick PCR Purification columns (Qiagen). Purified pooled DNA samples were eluted in a final volume of 10μl. DNA Fragments were amplified in 50 μl volume reactions containing 2 μl pooled DNA, 25 μl KAPA HiFi HotStart Master Mix (Kapa Biosystems), and 2 μl of both PCR primers (12.5 μM). PCR cycling consisted of 98° for 30 sec, followed by 18 cycles of 98° for 30 sec, 65° for 30 sec, 72° for 30 sec with a final extension for 5 min and kept at 4°. Amplified libraries were purified as above but eluted in 30 μl. Of the amplified libraries 1 μl was loaded onto a Bioanalyzer High Sensitivity DNA Chip (Agilent technologies) for evaluation of fragment sizes and 1 μl was used for quantification using Qubit (Life Technologies). Amplified library products were used for extra size selection using 2% agarose

gel cassette on a blue pippin system (Sage Science) to remove fragments smaller than 300 bp. Eluted size selected libraries were purified by AmpureXP beads (Agencourt). Final libraries were used for clustering on five lanes of an illumina Paired End flowcell using a cBot. Sequencing using an illumina HiSeq2000 instrument using 2\*100 nt Paired End reads.

### Sequencing and processing raw GBS data

Raw sequencing data were processed using the TASSEL software package (Glaubitz *et al.* 2014) version 5.2.37 using the GBSv2 pipeline. For quality filtering and barcode trimming, the GBSSeqToTagDBPlugin was run with the following parameters: kmerLength: 64, minKmerL: 20, mnQs: 20, mxKmerNum 100000000. Tags were dumped from the produced database using TagExportToFastqPlugin and mapped to the reference genome using the bwa software package (Li and Durbin 2010) in single-ended mode (samse). Positional information from aligned SAM files was stored in the TASSEL database using the SAMToGBSdbPlugin. The DiscoverySNPCallerPlugin was run using the following parameters: mnLCov: 0.1, mnMAF: 0.01. Found SNPs were scored for quality using SNPQualityProfilerPlugin and the Average taxon read depth at SNP was used as a quality score for filtering in the next step (minPosQS parameter), these scores were written to the TASSEL database using UpdateSNPPositionQualityPlugin. Finally, the ProductionSNPCallerPluginV2 was run with the following parameters: Avg Seq Error Rate: 0.002, minPosQS: 10, mnQS: 20.

### Genetic map calculation

We used JoinMap v4.1 for the genetic map construction (Stam 1993; van Ooijen 2006). The genetic map v2.5 was built with 749 SNPs generated by GBS based on genome v2.5 (unprocessed and processed data available as S1 and S2). A set of 632 SNPs called according to genome 3.0 was used for the genetic map v3.0 (unprocessed and processed data available as S3 and S4). Regression and Maximum likelihood mapping were used to calculate these maps (the linkage group information for both 2.5 and 3.0 genetic maps are available as S5).

**Genome improvement:** *Genome version vAM:* AllMaps We ran AllMaps (Tang *et al.* 2015) with default setting to combine genetic map v2.5 and physical map genome v2.5. This step resulted in genome vAM, in which scaffolds were ordered and oriented to reconstruct chromosomes.

### CONTAMINATION REMOVAL

The *Ae. arabicum* scaffolds v2.5 were checked for contaminations. The genome scaffolds were split into 197,702 1 kbp fragments and blasted against the NCBI nt database (NCBI Resource Coordinators 2016) using Tera-BLAST (TimeLogic Tera-BLAST algorithm, Active Motif Inc., Carlsbad, CA). The output was then analyzed by MEGAN 6 (Huson *et al.* 2016). All scaffolds for which more than 50% of their entire length was found in bacteria and with no hit in Viridiplantae were marked as putatively contaminated. A hit was counted with a minimum bit score of 50. Additionally, Bisulfite sequencing (BS-seq) CpG and Chromatin ImmunoPrecipitation DNA-Sequencing (ChIP-seq) H3 data were checked to identify contaminants (Aethionema\_contamination.xlsx). We ensured that these scaffolds were not combined with another scaffold by PBjelly2. After screening, three v2.5 scaffolds were removed: Scaffold\_2406, Scaffold\_2454 and Scaffold\_2594. They had a total length of 1,758 bp without any annotated genes. A summary of the contamination screen is available as S6.

■ **Table 1 (Mérat *et al.* 2019): Overview of the *Aethionema arabicum* PacBio reads**

	Total reads	Cyprus	Turkey
Number of reads	381,069	152,415	228,654
Length variation	11 - 57,910	11 - 55,919	11 - 57,910
Average length	5,845	5,795	5,879
Average quality	10.5	10.1	10.7

The lengths are given in nucleotides and the quality as phred score.

## LONG READ GENERATION FOR GENOME IMPROVEMENT

### PacBio reads

Genomic DNA (gDNA) for *Ae. arabicum* was obtained from leaves of the Cyprus and Turkey ecotypes. DNA was extracted using a modified protocol (Hiss *et al.* 2017) based on (Dellaporta *et al.* 1983). For the Turkey ecotype 35.70 µg and for the Cyprus ecotype 21.45 µg high molecular weight DNA were sent to the Max Planck-Genome-Centre, Cologne, and sequenced using the PacBio RS II machine (library insert size was 10-15 kbp gDNA). Four flow cells for Cyprus and six for Turkey were sequenced. Table 1 summarizes the statistics of the reads. The CG content of the pooled reads was 38%.

### MinION reads

To obtain MinION reads, gDNA was extracted from the Turkey ecotype (four leaf samples, 100 mg each) as outlined above. After pooling the samples, the gDNA concentration was measured using Hoechst 33258 DNA dye and resulted in 73.85 ng/µl. The library preparation was done using the Oxford Nanopore SQK-NSK007 protocol and R9.4 chemistry to design an 8 kbp 2D library. The sequencing run was carried out using Oxford Nanopores MinION technology. The flow cell sequenced 30,935 reads (122,362,072 nt) at -205 mV and 24 hr of runtime. After base calling with the MinKNOW 1.6 software (Oxford Nanopore Technologies Ltd.) the read length ranged from 5 to 63,441 nt with an average length of 3,955 nt. The average phred quality score was 11 and the GC content 41%, reads were not filtered or trimmed. The initial sequence format FAST5 was converted to FASTQ format by using the R package poRe version 0.21 (Watson *et al.* 2015). Because the MinION flow cell had previously been used for *Physcomitrella patens* DNA in the same run, the 30,935 reads were filtered for putative *P. patens* contamination. The reads were mapped with the long read mapper GMAP version 2017-08-15 (Wu and Nacu 2010) against the *P. patens* genome V3 (Lang *et al.* 2018). All settings were kept at default. 1,447 reads were characterized as putative *P. patens* reads and therefore removed.

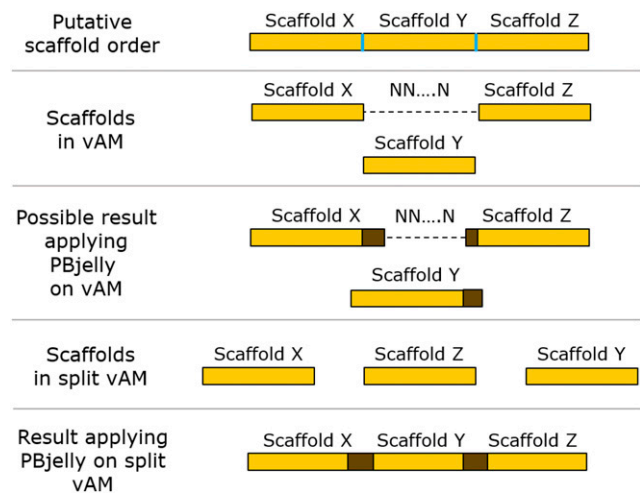
### Genome improvement using long reads

To perform super-scaffolding and gap filling, the program PBjelly 2 version 15.8.24 was used (English *et al.* 2012). It internally uses BLASR v5.1 (Chaisson and Tesler 2012) for mapping reads to a reference. PBjelly setup was used with a minimal gap size (-minGap) of 10 and called BLASR with -minMatch 12 (based on the observation that lower values yield less mappings), -bestn 1 (mapping reads uniquely) and -noSplitSubreads.

### Genome version 2.6: PacBio sequencing incorporation

We ran PBjelly2 with 381,022 (152,398 CYP, 228,624 TUR) PacBio reads which were head-cropped with 20 (due to suspicious per base sequence content suggestion presence of adapters) using Trimmomatic version 0.36 (Bolger *et al.* 2014).

PBjelly2 was used to improve genome v2.5 and vAM. Comparing the results, we found some scaffold connections which were made by



**Figure 2** Problem arising from applying PBjelly2 on vAM. Scaffold borders are visualized in blue and extensions of scaffolds introduced by PBjelly2 are shown in brown. Assuming the true order of the scaffolds is shown on top of the figure, but scaffold X and scaffold Z were already combined in the vAM assembly (second bar from top) this could lead to a partial filling of the N-stretch and maybe an extension of scaffold Y. However, PBjelly2 would not be able to place scaffold Y between the two other scaffolds (middle bar). If the scaffolds were thus split again (second bar from bottom), it is possible that the connections are made correctly applying PBjelly2 on the split version (bottom bar). This only visualizes a theoretical case, in this work it appeared every time that scaffold X and Y were connected by PBjelly2 and scaffold Z had to be reconnected afterward.

PBjelly2 (v2.5) were no longer possible for vAM (these scaffolds were already connected). Five connections formed for v2.5 scaffolds were already introduced by the genetic map approach (see above). Twelve connections which could be established in v2.5 were not formed in the PBjelly2 output for improving vAM, because the scaffolds were already connected with other scaffolds. Since PBjelly2 only fills gaps with reads and is not able to place whole scaffolds in gaps, it was necessary to split the vAM genome at certain points to be able to obtain the twelve connections which were not present in the PBjelly2 output for vAM (visualized in Figure 2). Split scaffolds were reconnected again after running PBjelly2, using N-stretches of length 100 to keep all improvements introduced in vAM if they were not formed by PBjelly2 (scaffolds in the vAM genome were combined using stretches of 100 Ns to denote a gap of unknown length). Since it is possible that PBjelly2 only fills a gap partially, we had to identify the positions of the gaps introduced by AllMaps in the new genome version and checked if they were filled completely or not. If the gap length was reduced, it was extended to have a length of 100 again. This approach produced genome v2.6.

### Genome version 3.0: MinION sequencing incorporation

After improving the genome to v2.6 using the PacBio reads, the same approach was applied for 30,935 MinION TUR reads to obtain the *Ae. arabicum* genome v3.0. The MinION reads were also checked for contamination. The genome version 3.0 is available at <https://genomevolution.org/coge/GenomeView.pl?gid=36061>.

### Name convention of *Ae. arabicum* v3.0 genome scaffolds

Scaffolds of genome v3.0 were named and ordered according to their length from long to short. The longest eleven scaffolds were named



**Table 2 Overview of gene liftover: GFF migration statistics**

Lifted only by flo	4,346
Lifted only by GeMoMa	36
Lifted with both programs	18,177
Manually lifted	34
Partially lifted	14
Number of corrected CDS	10,259
Marked as pseudo	3,230

linkage group (LG) based on the genetic map. Scaffolds which were combined are named csc for concatenated scaffold and the other ones are named sc (scaffold). The v3.0 scaffold names therefore follow the scheme type-number□v2.5 scaffold[v2.5 scaffold. . .]. I.e., the scaffold type (LG, CSC, SC), followed by a minus and the number of the scaffold, separated by a blank, followed a list of scaffolds denoting the v2.5 scaffolds or the v3.0 scaffolds. This naming system resulted in a shift in LG order between v2.5 and v3.0 (Supplementary file linkage\_group\_map.xlsx).

### Migration of proteins to new genome version

To perform the lift over of the gene models from v2.5 to v3.0, a combination of Gene Model Mapper (GeMoMa) v1.4 (Keilwagen *et al.* 2016) and flo (flo - same species annotations lift over pipeline, <https://github.com/wurmlab/flo>) were used. The results of both programs were concatenated. flo results were preferred over GeMoMa results if the results of the two programs differed, because flo works with alignments on nucleotide level while GeMoMa works with blasting proteins on amino acid level. If a protein could not be lifted completely, it is marked as partial in the resulting GFF (v3.0). A total of 34 genes had to be lifted manually, because they were either not lifted at all or only partially. If an intron could not be lifted, it was added by hand. If an exon or CDS could not be lifted, the new location was deriving from neighboring features which could be migrated to the new genome version. The location was then used to extract the nucleotide sequence from the genome using samtools v1.4 (Li *et al.* 2009). Only if the sequence was identical to the original sequence extracted from v2.5, the feature was migrated. This check was performed with ClustalW v2.1 (Larkin *et al.* 2007). After the migration step, the GFF file was checked and corrected. Genes which did not contain a start or a stop, contained internal stops or whose CDS sequence had a length not dividable by 3 were marked as potential pseudogenes with “pseudo=true”. To check if a gene contains internal stops each of its CDS features was checked individually for having at least one frame which results in no stop codons. Genes which were identical to other genes (start and stop position are equal) or were contained in other genes were removed. If the 3' CDS of a gene did not contain a stop codon but could be added by extending the CDS by three nucleotides, the CDS was corrected. The lifted genes were classified as shown in Table 2.

Most genes could be lifted by flo and GeMoMa. The reason why flo was able to lift more proteins is that GeMoMa works with protein sequences and the program was not able to generate proteins for 20,056 CDS features, either because a gene did not possess a CDS or because of faulty CDS sequences.

### Name convention of v3.0 genes

Old gene IDs were kept in the note attribute of the genes in the GFF and the linkage group numbers of the genetic map are also noted. The names of the genes were changed into Aa3typeNumberGenenumber: Aa for *Aethionema arabicum*, indicator genome version 3, followed by the type of scaffold, its number and the number of the gene (starting with 1 at the 5' end), e.g., Aa3LG1G2 or Aa3SC2601G1). For transcript isoforms (splice variants) this locus nomenclature can be extended by the number of the isoform (.X). Version 3.0 of the genome and all gene models are available at <https://genomeevolution.org/coge/GenomeView.pl?gid=36061>.

### Data availability

The genome version 2.5 is available at: <https://genomeevolution.org/coge/GenomeInfo.pl?gid=23428>. The genome version 3.0 is available at: <https://genomeevolution.org/coge/GenomeView.pl?gid=36061>. The GBS unprocessed and processed reads for genome mapping 2.5 and 3.0 are available as supplemental files S1-S4. The linkage group information for both 2.5 and 3.0 genetic maps are available as S5. A summary of the contamination screen is available as S6. MinION and PacBio single end reads are available from NCBI SRA (BioProject PRJNA558876). Supplemental material available at FigShare: <https://doi.org/10.25387/g3.8233055>.

## RESULTS AND DISCUSSION

### Reference genome improvement

The published draft version of the *Ae. arabicum* genome utilized the Ray assembler and contained 59,101 scaffolds with an N50 of 115,195bp (Haudry *et al.* 2013). Reassembly using the the AllPathsLG assembler and gap-closing using the SOAPdenovo GapCloser tool were used as a starting point for super-scaffolding. This resulted in a reassembly with 3,166 scaffolds, and a scaffold N<sub>50</sub> of 564,741bp labeled and published as version 2.5 on <https://genomeevolution.org/coge/>. The subsequent genome versions (vAM, v2.6 and v3.0) were obtained using linkage map and long read correction. The quality improvement of the genome is presented as the increase in total number of bases, reduced number of scaffolds and number of gaps, as well as bigger N<sub>50</sub> and smaller L<sub>50</sub> parameters (Table 3). In comparison with the starting genome v2.5, the final genome v3.0 has 9% less scaffolds (from 3,166 to 2,883). The overall length of genome v3.0 was extended from 196,005,095 to 203,449,326 bases (17% more) and the number of uncalled bases was reduced from 25,768,296 to 13,790,434 (from 13.2 to 6.8%) (Table 3).

**Table 3 Statistic overview of *Aethionema arabicum* genome versions**

Genome version	Draft	v2.5	vAM <sup>a</sup>	v2.6	v3.0
# Bases		196,005,095	196,022,695	203,150,143	203,449,326
# Scaffolds	59,101	3,166	2,990	2,895	2,883
# Scaffolds containing Ns		1,910	1,734	1,542	1,539
# Ns		25,768,296	25,785,896	13,946,922	13,790,434
N50	115,195	564,741	10,141,718	10,328,388	10,328,388
L50		56	9	9	9

<sup>a</sup>The vAM assembly includes all scaffolds; a total of 199 of the 3,166 v2.5 scaffolds were scaffolded via the genetic map into the 11 Linkage Groups (LG) of vAM (included in the 2,990 scaffolds), the 11 LG comprise 125,484,166 bp (64% of vAM).

■ **Table 4 Mapping efficiency of PBjelly2's mapping step.** The percentages in brackets give the percentage of the total number of reads (CYP, TUR or CYP + TUR). The line “# covered scaffolds” gives the number of scaffolds in which at least one read was mapped. Here, the number in brackets gives the percentage of the total number of scaffolds

Setup	PacBio vs. v2.5	PacBio vs. vAM	MinION vs. v2.5	MinION vs. v2.6
# mapped TUR reads	198,675 (86.9%)	198,629 (86.9%)	14,098 (45.6%)	15,886 (51.4%)
# mapped CYP reads	131,976 (86.6%)	131,942 (86.6%)	—	—
# total reads mapped	330,651 (86.8%)	330,571 (86.8%)	14,098 (45.6%)	15,886 (51.4%)
# unmapped reads	50,371 (13.2%)	50,451 (13.2%)	16,837 (54.4%)	15,049 (48.6%)
# scaffolds input genome	3,166	2,990	3,166	2,895
# covered scaffolds	2,971 (93.8%)	2,804 (93.8%)	1,689 (53.3%)	1,429 (49.4%)

## Genome improvement using long reads

**Read mapping efficiency:** The results of the mapping of the reads to the genome using PBjelly2 are summarized in Table 4. Almost the same number of PBjelly2 reads were mapped to genome v2.5 and vAM. However, it was important to apply PBjelly2 on both genomes in order to find scaffold connections which were not possible due to a combination of certain scaffolds in vAM (see supplementary file combination\_comparison\_pbjelly\_for\_v2.5\_vs\_pbjelly\_for\_vGM.ods for details). The genome v2.6 resulted by improving the split vAM genome using PBjelly2, reconnecting scaffolds and resizing gaps if needed. We also compared the results for improving v2.5 with vAM, but there were no new scaffold connections which were missed by improving the v2.6 version, so we did not perform a split step for improving the genome using the MinION reads. The mapping efficiency for the MinION reads is lower than for the PacBio reads due to a contamination of the reads (see Methods for details). There are 5.9% more MinION reads which were mapped to v2.6 than v2.5, demonstrating that the changes done to the genome are supported by the very long reads.

**The effect of PBjelly2 runs applied to the different assembly versions:** In comparing the improvement approaches for v2.5, the genetic map approach was able to combine the highest number of scaffolds, resulting in the  $L_{50}$  value lowered from 56 to 9 and the  $N_{50}$  value increased 20-fold. The increase of the  $N_{50}$  value in case of the PBjelly2 result (using the PacBio reads for improving scaffolds) results mainly from improvements of the shorter scaffolds. Comparing the PBjelly2 result for applying the PacBio reads to v2.5 and vAM shows that the reduction of scaffold number and increase of number of bases in the genome is similar (Table 5). MinION reads could also be used for v2.5 assembly improvement, however the results were not as good as for using PacBio reads, due to a much smaller number of reads. Improving the genome v2.6 with the MinION reads is also possible, but the improvement is not as good as for v2.5. This demonstrates that connections done using the PacBio reads are also supported by MinION reads.

Comparison of values for the different genome versions with the values for the PBjelly2 output is shown. The PBjelly2 outputs are denoted in the form “read type” vs. “genome version” to show which reads were used to improve which version of the genome. The result for PacBio vs. vAM was the basis for v2.6 and MinION vs. v2.6 was the basis for v3.0.

While PBjelly2 does not do as good a job as the genetic map approach in connecting scaffolds, its power is revealed by the gap filling. In genome v2.5 a total of 1,910 scaffolds contained uncalled bases. This number was reduced to 1,711 scaffolds (by 7.3%; Table 6) using PBjelly2 with PacBio reads. The exact number of uncalled bases in the v2.5 *A. arabicum* genome was 25,768,296 (Table 6). In the PBjelly2 result only 13,940,203 Ns (Table 6) are left, a reduction by 45.9%. Comparing this with the PBjelly2 result for the improvement of the vAM genome using PacBio reads (Table 6), more gaps were removed from the connected

genome. The number of scaffolds containing Ns was reduced by 10.8% and the overall number of Ns was reduced by 50.0%, while the overall percentage of Ns in the genome remained the same in the two results. Due to the small number of MinION reads, the improvement of the assembly is less pronounced.

This table gives an overview of the number of Ns in the different genome versions and the PBjelly2 results. For the number of scaffolds containing Ns, the percentage is given relative to the total number of scaffolds is given in brackets. For the number of Ns, the percentage is relative to the total number of bases in the respective assembly.

## Migration of proteins to new genome version

The genome v2.5 harbors 23,594 annotated protein coding genes. Eight of them could not be lifted because they were located next to a gap in the genome. Since it is possible that PBjelly2 changes the sequences around gaps, the sequences of the genes were not identical anymore and the programs were therefore not able to migrate some genes from one assembly version to another. We checked the expression of the genes which could not be lifted using Illumina RNA-seq data representing several developmental stages (data not shown) and found that all of them had almost no expression, as a result they were not lifted manually. In addition to some unlifted genes, there were 17 genes which could be lifted only partially due to the same reason. All the other 23,569 genes could be lifted. A set of 579 genes were removed due to being identical with other genes, and 140 genes were removed because they were located completely in another gene. A total of 1,202 genes have no starting methionine, 2,055 have no stop, 132 genes contain internal stops and for 1,019

■ **Table 5 Overview of the PBjelly2 result statistics for the different setups**

Setup	v2.5 assembly	PacBio vs. v2.5 (PBjelly super-scaffolding)	MinION vs. v2.5 (PBjelly super-scaffolding)
# scaffolds	3,166	3,066	3,123
# bases	196,005,095	203,024,676	196,600,700
N50	564,741	542,490	564,741
L50	56	58	56
Setup	vAM	PacBio vs. vAM	
# scaffolds	2,990	2,905	
# bases	196,022,695	203,137,854	
N50	10,141,718	10,314,234	
L50	9	9	
Setup	v2.6	MinION vs. v2.6	
# scaffolds	2,895	2,886	
# bases	203,150,143	203,450,934	
N50	10,328,388	10,323,234	
L50	9	9	

Table 6 Gap/N analysis of different genome versions

Setup	v2.5	PacBio vs. v2.5	MinION vs. v2.5
# scaffolds containing Ns	1,910 (60.3%)	1,711 (56.0%)	1,901 (60.0%)
# Ns	25,768,296 (13.2%)	13,940,203 (7.1%)	25,142,571 (12.8%)
Setup	vAM	PacBio vs. vAM	
# scaffolds containing Ns	1,734 (58.0%)	1,546 (51.7%)	
# Ns	25,785,896 (13.2%)	13,942,094 (7.1%)	
Setup	v2.6		MinION vs. v2.6
# scaffolds containing Ns	1,542 (53.3%)		1,539 (53.2%)
# Ns	13,946,922 (6.9%)		13,790,284 (6.8%)

genes the length of the CDSs is not dividable by three. In the end 19,363 genes were lifted which were not marked as potential pseudogenes or partial.

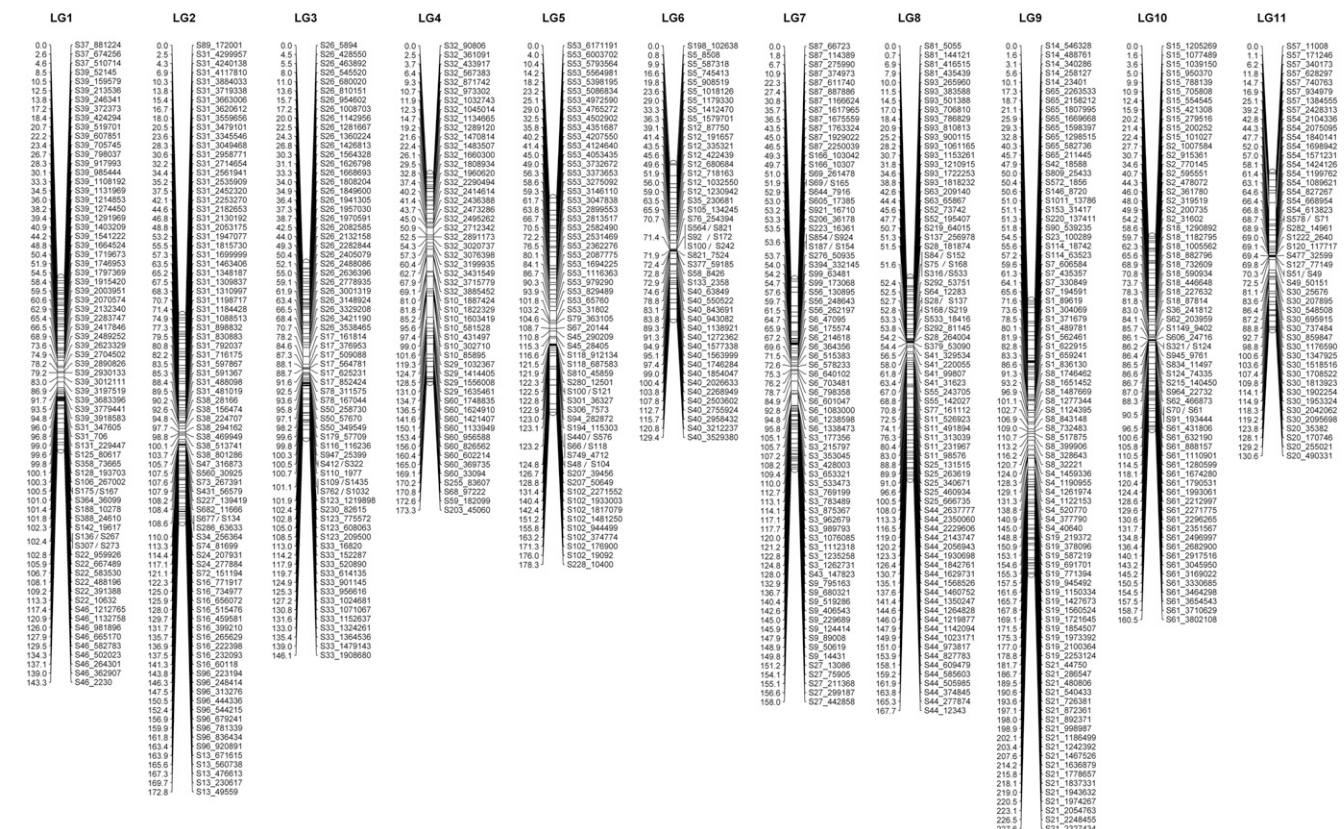
We find that the starting point (assembly version) for improvement is not relevant. PBjelly2 is able to make more improvements using the PacBio data than with the MinION data due to the higher number of PacBio reads. The number of added bases per read is higher for PacBio than MinION reads (18.71 vs. 9.67) and also the number of removed Ns is higher (31.07 vs. 5.06) while the MinION reads lead to more scaffold connections per read ( $2.49 \times 10^{-4}$  vs.  $3.88 \times 10^{-4}$  connected scaffolds). Using the MinION reads for improvement makes only a few changes, but they show that they support the changes which were made to the genome using the PacBio data. Since the *Ae. arabicum* genome was almost not contaminated at all, only three small scaffolds had to be removed. Gene models were filtered for multiple genes and genes

contained in other genes. If a problem with a gene was found, it was marked in the resulting GFF file. We note that there are gene models which are probably not correct and need to be fixed in future studies.

The combination of genetic mapping and long reads significantly improved the structure of the assembly, reducing the total scaffold number and decreasing the number of gaps.

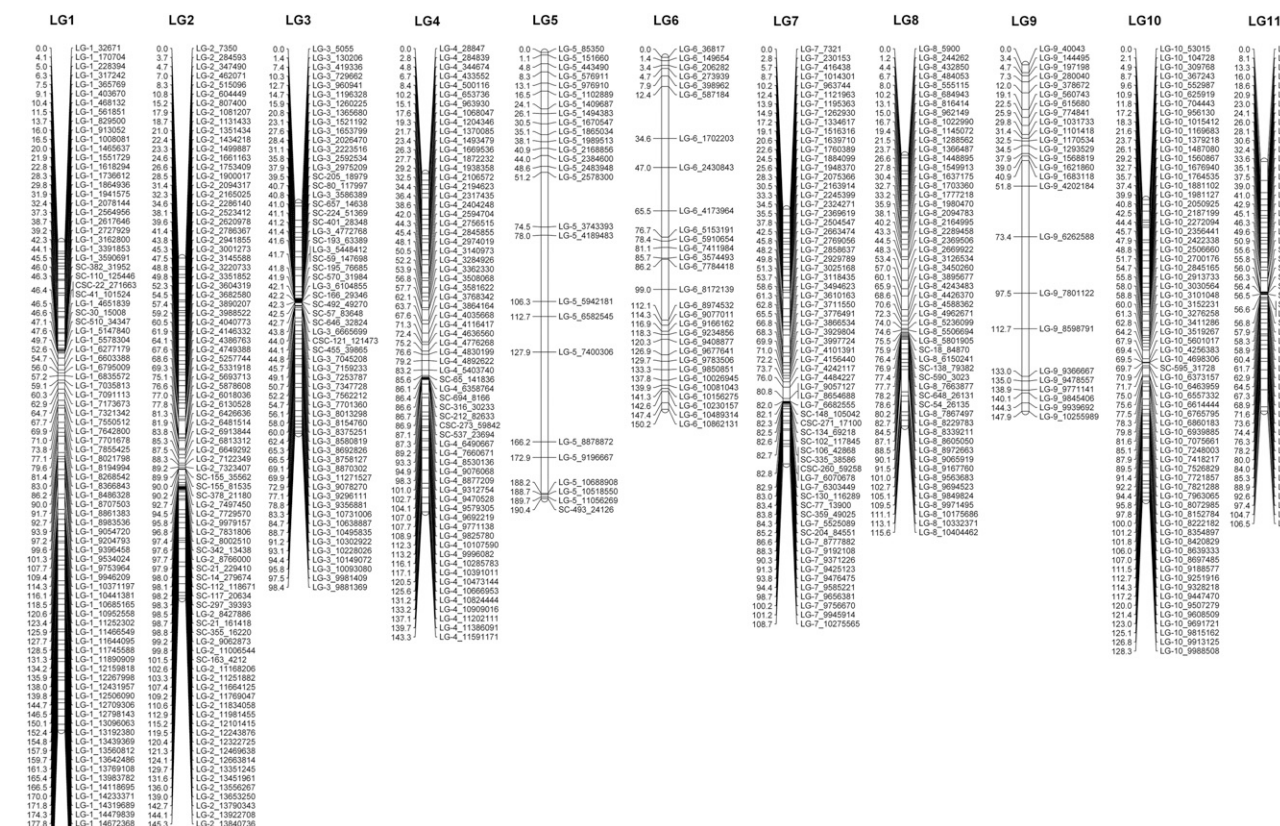
### Genetic map of *Aethionema arabicum*

**Genetic map v2.5: SNP calling** A GBS approach (Rowan *et al.* 2015) was used to generate genetic variation data for genetic mapping. Illumina sequencing of the parental lines and the RILs resulted in 442,101,405 raw reads after quality filtering. Using the TASSEL package (Glaubitz *et al.* 2014) to match sequence tags to markers, 160,379 SNPs could be called based on genome v2.5. SNPs identified through the GBS



**Figure 3** *Aethionema arabicum* genetic map v2.5. Genetic map version 2.5 consists of eleven linkage groups. On each linkage group, genetic distance in cM is present on the left and SNP markers on the right.





**Figure 4** *Aethionema arabicum* genetic map v3.0. Genetic map version 3.0 consists of eleven linkage groups. On each linkage group, genetic distance in cM is present on the left and SNP markers on the right.

method often take the form of many SNP ‘islands’, where a multitude of SNPs are present over only a few kbp of sequence with the same states across individuals. This makes genetic mapping difficult as it results in a very large number of markers that are mostly redundant. This is often referred to as the “large  $p$ , small  $n$ ” problem where  $p$  is the number of markers and  $n$  is the number of individuals (or recombination events) being used for mapping (Ronin *et al.* 2010). This leads to a situation where the number of markers greatly exceeds the resolution of recombination for the population used; thus, only a fraction of the markers can accurately be ordered. Another factor is that many markers cluster (co-segregate) to the same position in the map or are on the same genomic scaffold for example in our study. When there are many closely-linked markers, then any amount of genotyping errors leads to compounding and increasing estimates of recombination (map inflation). We reduced this SNP amount using a sliding window approach collapsing a group of SNPs that all have the same states across individuals into one single marker over windows of 10 kbp, thus the bigger the scaffold the more selected SNPs. This, together with filtering non-informative SNPs (missing data in more than 30% many individuals) resulted in a core group of markers of 5,428 SNPs.

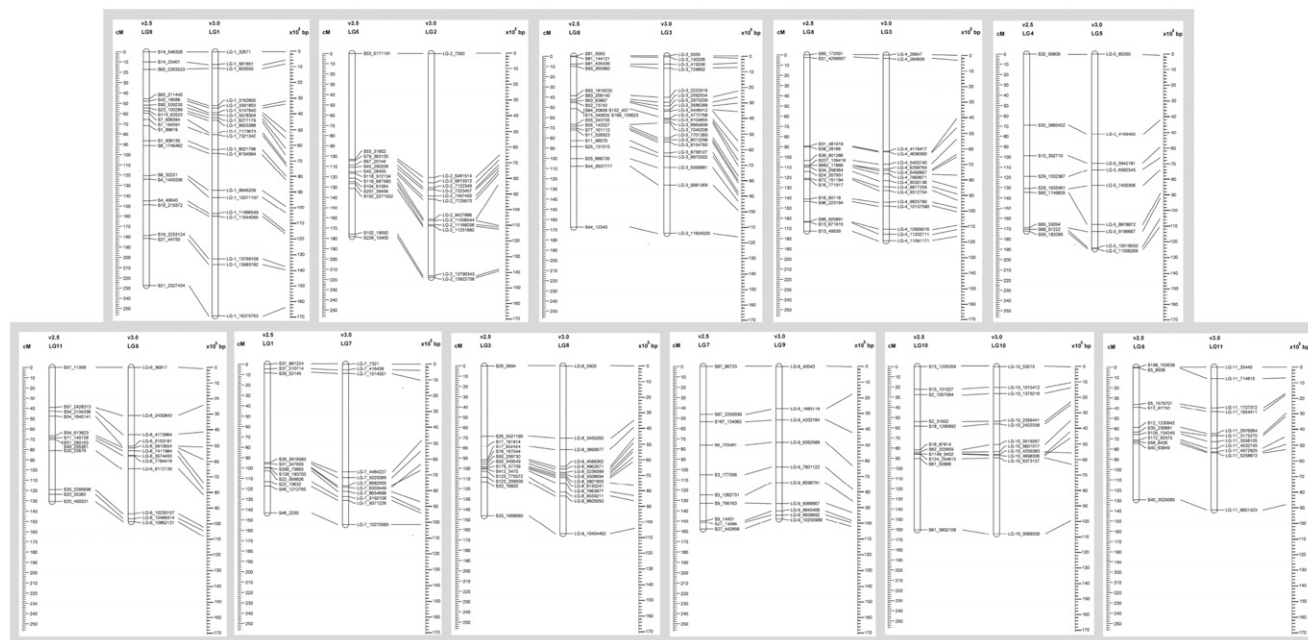
**Genetic map calculation** We used JoinMap 4.1 to calculate the genetic linkage map for the *Ae. arabicum* RIL population. For map v2.5, we first checked the marker similarity among the initial set of 5,428 SNPs that were obtained from GBS based on genome v2.5 by a pairwise

comparison. SNPs that were highly similar (higher than 90%) were represented by one marker, which refined the number of markers to 1,818. Grouping was selected at a LOD threshold of 9.0, which led to the grouping of the expected 11 Linkage Groups (LGs) (Figure 3). We further optimized each LG to avoid inflation of the map distance due to saturating SNPs using a Maximum likelihood model. Markers that were not more than 0.1 cM away were all eliminated. As a result, a final set of 746 SNP markers was used on 11LGs. Out of 11 LGs, there are three LGs (4, 7 and 11) containing cluster of segregation distorted SNPs (more than 50% of SNPs, supplemental file linkage\_group\_map.xlsx).

The *Ae. arabicum* genetic map v2.5 consists of 11 LGs with average distance size of 162.5 cM, covered by 746 SNP markers with average of 67 markers per LG. The average marker spacing was 2.4 cM, equivalent to approximately 169 kbp. The centromere is suggested by the high density of SNP markers within a small genetic distance (e.g., a low recombination frequency). These markers typically belong to relatively small scaffolds, consistent with a high-repetitive DNA content, where only a few SNPs were called. LG4 centromere is located at the end of the linkage, suggesting an assembly problem or that LG4 is a telocentric chromosome.

Overall the markers are distributed relatively dense and even in v2.5, with the biggest gap smaller than 18 cM. SNPs that reside in the same scaffold were in agreement among each other on the direction of their scaffold.





**Figure 5** The alignment of genetic map v2.5, v3.0 and physical map. The alignment of the genetic map v2.5 and v3.0 were based on relative SNPs. The left ruler indicates genetic distance in cM and the right indicates physical distance in bp according to genome v3.0.

**Genetic map v3.0:** The procedure to build genetic map v3.0 was similar to v2.5. SNP calling was performed based on genome v3.0 instead of v2.5 resulting in a raw set of 141,914 SNPs. After similar quality control strategy as for v2.5, we construct v3.0 with a core set of 632 SNPs (Figure 4). The 11 LGs were maintained with the total size of ~1945 cM, average marker distance of 3.1 cM. This inflation of genetic distance in genetic map v3.0 compared with v2.5 can be explain by the newly retrieved SNPs due to resolved Ns in the genome. These new SNPs are mainly located in the centromeric regions. In general, SNP order and orientation in LGs are in agreement with map v2.5 (Figure 5). We have made adjusted linkage group order in map v3.0 compared with map v2.5 (linkage\_group\_map.xlsx). Based on the size of the group, the biggest one is LG 1, and the smallest is LG 11 (Figure 4).

However, there is a significant difference between v3.0 and v2.5 at three LGs that harbor clusters of segregation distorted SNPs, LG 5, 6, and 9 (equivalent LG 4, 11 and 7 in v2.5, respectively): the reduced number of markers as well as the increased marker distance (Figure 3 and 4). In order to maintain 11 LGs and certain degree of newly called SNP incorporation, we had to reduced number of distorted markers in those LGs in v3.0, as a result the dearth of markers was observed (see supplemental file. linkage\_group\_map.xlsx).

## CONCLUSIONS AND POTENTIAL IMPLICATIONS

*Aethionema* is becoming an outgroup model for the rest of the core Brassicaceae. Studies on its genome, relevant life-history traits and their evolution rely on high-quality genomic and genetic resources. Thus, this work helps pave the way for future comparative genomic and trait research of the entire Brassicaceae family. We have provided an advanced version of the *Aethionema arabicum* genome and its first genetic map, which allows for pseudochromosome construction needed for analysis of genome evolution. It should be noted that due to the liftover procedure no previously undiscovered gene models were added; future work will need to add and refine gene models. Finally, quantitative trait loci (QTL) mapping for the wide range of traits in *Aethionema* (e.g., flowering time, chemical defenses (Hofberger *et al.* 2013;

Mohammadin *et al.* 2018), fruit heteromorphism (Lenser *et al.* 2016; Wilhelmsson *et al.* 2019; Arshad *et al.* 2019) and light-control of germination (Mérat *et al.* 2019) will be greatly enabled by the genetic map and improved genome assembly.

## ACKNOWLEDGMENTS

We thank Dr. Elio Schijlen for technical support on Genotyping by Sequencing technology, Dr. Lidija Berke for her help on running AllMaps and Marco Göttig, Rabea Meyberg, Christopher Grosche and Manuel Hiss for their help with long read sequencing. We also thank the other members of the SeedAdapt Consortium and Dr. Laurie Grandont for fruitful discussions about the work.

## LITERATURE CITED

- Al-Shehbaz, I. A., M. A. Beilstein, and E. A. Kellogg, 2006 Systematics and phylogeny of the Brassicaceae (Cruciferae): an overview. *Plant Syst. Evol.* 259: 89–120. <https://doi.org/10.1007/s00606-006-0415-z>
- Arshad, W., F. Marone, M. E. Collinson, G. Leubner-Metzger, and T. Steinbrecher, 2019 Fracture of the dimorphic fruits of *Aethionema arabicum* (Brassicaceae). *Botany* 1–11. <https://doi.org/10.1139/cjb-2019-0014>
- Beilstein, M. A., I. A. Al-Shehbaz, S. Mathews, and E. A. Kellogg, 2008 Brassicaceae phylogeny inferred from phytochrome A and ndhF sequence data: tribes and trichomes revisited. *Am. J. Bot.* 95: 1307–1327. <https://doi.org/10.3732/ajb.0800065>
- Bibalani, G. H., 2012 Investigation on flowering phenology of Brassicaceae in the Shanjian region Shabestar district, NW Iran (usage for honeybees).
- Boisvert, S., F. Laviolette, and J. Corbeil, 2010 Ray: Simultaneous Assembly of Reads from a Mix of High-Throughput Sequencing Technologies. *J. Comput. Biol.* 17: 1519–1533. <https://doi.org/10.1089/cmb.2009.0238>
- Bolger, A. M., M. Lohse, and B. Usadel, 2014 Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* 30: 2114–2120. <https://doi.org/10.1093/bioinformatics/btu170>
- Chaisson, M. J., and G. Tesler, 2012 Mapping single molecule sequencing reads using basic local alignment with successive refinement (BLASR):

- application and theory. *BMC Bioinformatics* 13: 238. <https://doi.org/10.1186/1471-2105-13-238>
- Dellaporta, S. L., J. Wood, and J. B. Hicks, 1983 A plant DNA miniprep: Version II. *Plant Mol. Biol. Report.* 1: 19–21. <https://doi.org/10.1007/BF02712670>
- Doyle, J., 1991 DNA Protocols for Plants, pp. 283–293 in *Molecular Techniques in Taxonomy*, edited by G. M. Hewitt, A. W. B. Johnston, and J. P. W. Young. NATO ASI Series, Springer, Berlin Heidelberg.
- Elshire, R. J., J. C. Glaubitz, Q. Sun, J. A. Poland, K. Kawamoto *et al.*, 2011 A Robust, Simple Genotyping-by-Sequencing (GBS) Approach for High Diversity Species. *PLoS One* 6: e19379. <https://doi.org/10.1371/journal.pone.0019379>
- English, A. C., S. Richards, Y. Han, M. Wang, V. Vee *et al.*, 2012 Mind the gap: upgrading genomes with Pacific Biosciences RS long-read sequencing technology. *PLoS One* 7: e47768. <https://doi.org/10.1371/journal.pone.0047768>
- Franzke, A., M. A. Lysak, I. A. Al-Shehbaz, M. A. Koch, and K. Mummenhoff, 2011 Cabbage family affairs: the evolutionary history of Brassicaceae. *Trends Plant Sci.* 16: 108–116. <https://doi.org/10.1016/j.tplants.2010.11.005>
- Glaubitz, J. C., T. M. Casstevens, F. Lu, J. Harriman, R. J. Elshire *et al.*, 2014 TASSEL-GBS: A High Capacity Genotyping by Sequencing Analysis Pipeline. *PLoS One* 9: e90346. <https://doi.org/10.1371/journal.pone.0090346>
- Gnerre, S., I. MacCallum, D. Przybylski, F. J. Ribeiro, J. N. Burton *et al.*, 2011 High-quality draft assemblies of mammalian genomes from massively parallel sequence data. *Proc. Natl. Acad. Sci. USA* 108: 1513–1518. <https://doi.org/10.1073/pnas.1017351108>
- Guo, X., J. Liu, G. Hao, L. Zhang, K. Mao *et al.*, 2017 Plastome phylogeny and early diversification of Brassicaceae. *BMC Genomics* 18: 176. <https://doi.org/10.1186/s12864-017-3555-3>
- Haudry, A., A. E. Platts, E. Vello, D. R. Hoen, M. Leclercq *et al.*, 2013 An atlas of over 90,000 conserved noncoding sequences provides insight into crucifer regulatory regions. *Nat. Genet.* 45: 891–898. <https://doi.org/10.1038/ng.2684>
- Hiss, M., R. Meyberg, J. Westermann, F. B. Haas, L. Schneider *et al.*, 2017 Sexual reproduction, sporophyte development and molecular variation in the model moss *Physcomitrella patens*: introducing the ecotype Reute. *Plant J.* 90: 606–620. <https://doi.org/10.1111/tpl.13501>
- Hofberger, J. A., E. Lyons, P. P. Edger, J. Chris Pires, and M. Eric Schranz, 2013 Whole genome and tandem duplicate retention facilitated glucosinolate pathway diversification in the mustard family. *Genome Biol. Evol.* 5: 2155–2173. <https://doi.org/10.1093/gbe/evt162>
- Huang, C.-H., R. Sun, Y. Hu, L. Zeng, N. Zhang *et al.*, 2016 Resolution of Brassicaceae Phylogeny Using Nuclear Genes Uncovers Nested Radiations and Supports Convergent Morphological Evolution. *Mol. Biol. Evol.* 33: 394–412. <https://doi.org/10.1093/molbev/msv226>
- Huson, D. H., S. Beier, I. Flade, A. Górski, M. El-Hadidi *et al.*, 2016 MEGAN Community Edition - Interactive Exploration and Analysis of Large-Scale Microbiome Sequencing Data. *PLOS Comput. Biol.* 12: e1004957. <https://doi.org/10.1371/journal.pcbi.1004957>
- Imbert, E., 2002 Ecological consequences and ontogeny of seed heteromorphism. *Perspect. Plant Ecol. Evol. Syst.* 5: 13–36. <https://doi.org/10.1078/1433-8319-00021>
- Keilwagen, J., M. Wenk, J. L. Erickson, M. H. Schattat, J. Grau *et al.*, 2016 Using intron position conservation for homology-based gene prediction. *Nucleic Acids Res.* 44: e89. <https://doi.org/10.1093/nar/gkw092>
- Kent, W. J., C. W. Sugnet, T. S. Furey, K. M. Roskin, T. H. Pringle *et al.*, 2002 The Human Genome Browser at UCSC. *Genome Res.* 12: 996–1006. <https://doi.org/10.1101/gr.229102>
- Lang, D., K. K. Ullrich, F. Murat, J. Fuchs, J. Jenkins *et al.*, 2018 The *Physcomitrella patens* chromosome-scale assembly reveals moss genome structure and evolution. *Plant J.* 93: 515–533. <https://doi.org/10.1111/tpl.13801>
- Larkin, M. A., G. Blackshields, N. P. Brown, R. Chenna, P. A. McGettigan *et al.*, 2007 Clustal W and Clustal X version 2.0. *Bioinformatics* 23: 2947–2948. <https://doi.org/10.1093/bioinformatics/btm404>
- Lenser, T., K. Graeber, Ö. S. Cevik, N. Adigüzel, A. A. Dönmez *et al.*, 2016 Developmental Control and Plasticity of Fruit and Seed Dimorphism in *Aethionema arabicum*. *Plant Physiol.* 172: 1691–1707. <https://doi.org/10.1104/pp.16.00838>
- Li, H., and R. Durbin, 2010 Fast and accurate long-read alignment with Burrows–Wheeler transform. *Bioinformatics* 26: 589–595. <https://doi.org/10.1093/bioinformatics/btp698>
- Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan *et al.*, 2009 The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25: 2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>
- Luo, R., B. Liu, Y. Xie, Z. Li, W. Huang *et al.*, 2012 SOAPdenovo2: an empirically improved memory-efficient short-read de novo assembler. *Gigascience* 1: 18. <https://doi.org/10.1186/2047-217X-1-18>
- Mérai, Z., K. Graeber, P. Wilhelmsson, K. K. Ullrich, W. Arshad *et al.*, 2019 *Aethionema arabicum*: a novel model plant to study the light control of seed germination. *J. Exp. Bot.* 70: 3313–3328. <https://doi.org/10.1093/jxb/erz146>
- Mohammadin, S., K. Peterse, S. J. van de Kerke, L. W. Chatrou, A. A. Dönmez *et al.*, 2017 Anatolian origins and diversification of *Aethionema*, the sister lineage of the core Brassicaceae. *Am. J. Bot.* 104: 1042–1054. <https://doi.org/10.3732/ajb.1700091>
- Mohammadin, S., W. Wang, T. Liu, H. Moazzeni, K. Ertugrul *et al.*, 2018 Genome-wide nucleotide diversity and associations with geography, ploidy level and glucosinolate profiles in *Aethionema arabicum* (Brassicaceae). *Plant Syst. Evol.* 304: 619–630. <https://doi.org/10.1007/s00606-018-1494-3>
- NCBI Resource Coordinators, 2016 Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res.* 44: D7–D19. <https://doi.org/10.1093/nar/gkv1290>
- Ronin, Y., D. Mester, D. Minkov, and A. Korol, 2010 Building reliable genetic maps: different mapping strategies may result in different maps. *Nat. Sci.* 02: 576.
- Rowan, B. A., V. Patel, D. Weigel, and K. Schneeberger, 2015 Rapid and Inexpensive Whole-Genome Genotyping-by-Sequencing for Crossover Localization and Fine-Scale Genetic Mapping. *G3: Genes, Genomes, Genetics* 5: 385–398.
- Stam, P., 1993 Construction of integrated genetic linkage maps by means of a new computer package: Join Map. *Plant J.* 3: 739–744. <https://doi.org/10.1111/j.1365-313X.1993.00739.x>
- Tang, H., X. Zhang, C. Miao, J. Zhang, R. Ming *et al.*, 2015 ALLMAPS: robust scaffold ordering based on multiple maps. *Genome Biol.* 16: 3. <https://doi.org/10.1186/s13059-014-0573-1>
- van Ooijen, J. W., 2006 JoinMap 4. Software for the calculation of genetic linkage maps in experimental populations.
- Watson, M., M. Thomson, J. Risse, R. Talbot, J. Santoyo-Lopez *et al.*, 2015 poRe: an R package for the visualization and analysis of nanopore sequencing data. *Bioinformatics* 31: 114–115. <https://doi.org/10.1093/bioinformatics/btu590>
- Wilhelmsson, P. K. I., J. O. Chandler, N. Fernandez-Pozo, K. Graeber, K. K. Ullrich *et al.*, 2019 Usability of reference-free transcriptome assemblies for detection of differential expression: a case study on *Aethionema arabicum* dimorphic seeds. *BMC Genomics* 20: 95. <https://doi.org/10.1186/s12864-019-5452-4>
- Wu, T. D., and S. Nacu, 2010 Fast and SNP-tolerant detection of complex variants and splicing in short reads. *Bioinformatics* 26: 873–881. <https://doi.org/10.1093/bioinformatics/btq057>

Communicating editor: J. Ross-Ibarra