



Identification of *Arabis alpina* genomic regions associated with climatic variables along an elevation gradient through whole genome scan

Stéphane Lobréaux, Christian Miquel

► To cite this version:

Stéphane Lobréaux, Christian Miquel. Identification of *Arabis alpina* genomic regions associated with climatic variables along an elevation gradient through whole genome scan. Genomics, Elsevier, 2019, 10.1016/j.ygeno.2019.05.008 . hal-02336683

HAL Id: hal-02336683

<https://hal.archives-ouvertes.fr/hal-02336683>

Submitted on 29 Oct 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 **Identification of *Arabis alpina* genomic regions associated with climatic**
2 **variables along an elevation gradient through whole genome scan**

3
4
5
6
7 Authors : Stéphane Lobréaux^{1*}, Christian Miquel¹

8
9
10
11 ¹ Laboratoire d'Ecologie Alpine, CNRS, Université Grenoble-Alpes, Grenoble, France

12
13 * corresponding author : stephane.lobreaux@univ-grenoble-alpes.fr

14
15
16
17 Keywords : SNP, local adaptation, plant, membrane-bound transcription factor, genome scan

18
19
20
21
22 Data accessibility : Sequence data from this article have been deposited with the European
23 Nucleotide Archive under the project number PRJEB32228.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Abstract

We performed a pooled whole-genome sequencing on samples of the alpine plant *Arabis alpina*, harvested in ten populations along an elevation gradient in the French Alps. A large dataset of genetic variations was produced as single nucleotide polymorphisms (SNPs). A combined genome scan approach enabled detecting genomic regions associated with a synthetic environmental variable characterizing the climate at each sampling location. Positive loci detected by two methods were retained and belong to 19 regions in the *Arabis alpina* genome. The most significant region harbors an ortholog of the AtNAC062 gene, encoding a membrane-bound transcription factor described as linking the cold response and pathogen resistance that may confer protection to plants under extended snow coverage at high elevations. Other genes involved in the stress response or in flowering regulation were also detected. Altogether, our results indicated that *Arabis alpina* represent a suitable model for studying genomic adaptation in alpine perennial plants.

1 **1. Introduction**

2 Local adaptation is a process by which a population adapts to its local biotic and/or
3 abiotic environmental conditions and becomes well adapted to this particular environment [1].
4 Local adaptation is an important field of research that has been stimulated by the recent
5 development of new methods and the interest in ongoing climatic change [2, 3]. At the genome
6 level, genetic variations related to local adaptation tend to be maintained into the population,
7 and the corresponding genomic regions are therefore subjected to selective pressure [4]. At
8 these loci, a change in allelic frequencies can be detected in response to environmental
9 constraints, and different methods based on statistical tools have been developed to look for
10 such genomic regions [5, 6]. In searching for genome-environment associations, it is important
11 to be aware of confounding parameters, such as population structure and demography, which
12 may lead to detecting false positives [4]. Evolution in sequencing technologies has enabled to
13 acquire genomic data and genetic variations at a larger scale more easily and at a reduced cost
14 [7]. Some techniques targeting genomic regions by reduced representations such as GBS [8] or
15 RADseq [9] allow further reducing genotyping cost. However, genome resequencing remains,
16 when affordable, a valuable strategy for obtaining SNP data at high density and along the entire
17 genome [10] to get a more accurate and exhaustive view of genetic variations associated with
18 environmental data. An alternative strategy for acquiring allelic frequency estimates from
19 populations is pooled resequencing of individuals [11-13]. This strategy makes it possible to
20 significantly reduce sequencing costs while SNP density and genome coverage are preserved.
21 Different studies have validated and successfully implemented this strategy in the past years to
22 investigate local adaptation in plants, for example in *Zea mays* [14] or *Populus* trees [15].

23

24 Mountain areas are suitable environments to look for local adaptation since contrasting
25 climatic conditions occur within a rather short distance from each other [16]. It is well described

1 that along elevation gradients, living organisms face harsher climatic conditions with an
2 increase of elevation. Studying populations of a species along such gradients may help to
3 identify genomic regions involved in the local adaptation of a particular species. Due to their
4 immobility, plants are particularly exposed to local climate conditions, and their adaptation is
5 therefore vital. Alpine plants face exposure to stress factors such as UV, cold, high light, wind,
6 and extended snow coverage [17]. Among alpine plants, some of them, like *Arabis alpina*, are
7 able to grow in a wide range of elevations. This plant of the Brassicaceae family is found in the
8 French Alps at elevations varying from 500 to 3,000 m. Such a large range of climatic
9 conditions makes it possible to study local adaptation by comparing genomic data from
10 populations sampled along such a climatic gradient. *Arabis alpina* has a relatively small
11 genome of 337 Mb for which data are available [18, 19], and the sequenced genome has recently
12 been improved using a combination of optical mapping and chromosome conformation capture
13 [20]. This plant has emerged as a model, closely related to *Arabidopsis thaliana*, allowing to
14 study the regulation of flowering in a perennial plant [21, 22]. But its capacity of growing in a
15 large range of elevations makes it also very interesting for studying local adaptation. Previous
16 studies have been conducted to detect genomic regions in *Arabis alpina* associated with local
17 environments [23, 24]. Poncet *et al.* [23] searched for ecologically relevant loci associated with
18 climatic variables based on samples harvested in the French and Swiss Alps. Both studies were
19 based on amplified fragment length polymorphisms (AFLP), and with this method, a limited
20 number of loci were investigated in the genome by a few hundred of AFLP markers.

21

22 In this study, we have investigated *Arabis alpina* local adaptation at the whole genome level by
23 resequencing pools of individuals sampled in populations along an altitude gradient in the
24 French Alps.

25

1 **1. Materials and methods**

2 *1.1 Sampling*

3 The ten sampling sites are located in the French Alps, covering a range of elevations from 860
4 to 2,781 m (see Table1 for sampling site data). Plant material was harvested from all sites during
5 a 15-day period in the summer. Fresh leaf samples were collected and dehydrated in silica gel,
6 in which they were stored until DNA extraction.

7

8 *1.2 DNA extraction and sequencing*

9 Genomic DNA was extracted from ten samples from each sampling location. For each sample,
10 twenty milligrams of dried leaves were grounded into a fine powder, and DNA extraction was
11 performed using the NucleoSpin Plant II kit (Macherey-Nagel, Duren, Germany) according to
12 the manufacturer's instructions. Extracted DNA was quantified using the Qubit dsDNA BR
13 assay kit on a Qubit fluorometer (Invitrogen, USA). Agarose gel electrophoresis enabled
14 checking DNA quality and verifying that the quantification was accurate. DNA pools were
15 prepared by mixing equal amounts of genomic DNA from the 10 samples of a location.

16 DNA sequencing using Illumina technology was performed by Fasteris (Geneva, Switzerland).
17 To prepare libraries from each DNA pool, DNA was broken, and bead size selection enabled
18 recovering inserts of 350bp. Libraries were prepared using the Illumina Genomic Nano kit.
19 Sequencing was performed using multiplexing on a HiSeq2500 sequencer producing reads of
20 150bp as paired-end data. Libraries were loaded on 3 flow cell lanes to get an average genome
21 coverage of 28x.

22

23 *2.3 SNP mapping*

24 Raw reads were filtered to remove sequences with an average quality below 25, and 5' and 3'
25 trimming was performed to delete bases with quality below 20. Any read presenting a length

1 lower than 50 bases after quality trimming was removed. Prior to read mapping on *Arabis*
2 *alpina* chromosomes, a filtering step to get rid of abundant chloroplast reads was performed.
3 Paired-end reads were mapped using BWA 0.7.14 [25] on *Arabis alpina* chloroplast genome
4 (accession number HF934132, [26]). Then, IDs from all matching reads were collected, and
5 the corresponding paired-end reads were removed from the dataset.

6 The resulting filtered dataset was mapped onto the *Arabis alpina* chromosomes. The
7 last published version 5 of this genome was used in this study [20]. Mapping was performed
8 using BWA 0.7.14, while samtools 1.3.1 and bcftools 1.3.1 enabled SNP calling and filtering
9 [27]. SNPs located within a distance of 10 bases around an indel were filtered from the SNP
10 dataset. We removed positions where more than 2 alleles were detected, where the coverage
11 was lower than 10, or where the mapping quality was below 30. We retained only SNPs for
12 which data were available in all 10 studied populations. When SNP variability was detected in
13 only one location, the corresponding SNP was discarded from the final dataset. A SNP density
14 filtering step was performed in which SNPs were removed when 3 SNPs or more were found
15 within a 10 base region.

16

17 *2.4 Genetic distance and phylogenetic tree*

18 One thousand SNPs were selected randomly in the mapping dataset from each of the 8
19 *Arabis alpina* chromosomes. Allelic frequencies of these 8,000 SNPs in the studied populations
20 were used to generate a phylogenetic tree by the neighbor-joining method using POPTREEW
21 [28]. The corresponding data in Newick format enabled plotting the tree with the ape R package
22 [29].

23

24 *2.5 Detection of loci associated with environmental variables*

25 *2.5.1 Environmental variables*

1 Climatic data were downloaded from the WorldClim database
2 (<http://www.worldclim.org>) at a spatial resolution of 30 arcsec. Precipitations and minimum
3 and maximum temperature data were extracted using the R package raster according to the GPS
4 coordinates of the sampling locations (Table 1). From these monthly data, annual precipitations
5 and the annual range of temperatures were calculated for each sampling site.

6
7 Snow cover data were downloaded from the National Snow and Ice Data Center
8 (<http://nsdic.org>). Such data are based on the difference in reflectance of snow-covered land,
9 which is high in the visible band but low in shortwave infrared. MODIS/Aqua Snow Cover 8-
10 Day L3 Global 500m Grid Version 6 files for the 2010-2016 period were used. This dataset
11 reports the maximum snow cover extent during an eight-day period in 1,200 km x 1,200 km
12 tiles. According to the GPS coordinates of the sampling sites, the appropriate tiles were defined
13 using the MODLAND tile calculator ([https://landweb.modaps.eosdis.nasa.gov/cgi-
14 bin/developer/tilemap.cgi](https://landweb.modaps.eosdis.nasa.gov/cgi-bin/developer/tilemap.cgi)). An average frequency of snow cover was calculated for each week
15 throughout the period 2010-2016. The number of weeks without snow was then defined,
16 corresponding to the sum of weeks where the snow cover frequency was below 0.5.

17
18 Principal component analysis (PCA) of environmental variables was performed with R
19 package ade4 [30].

20

21 *2.5.2 Generalized linear mixed models (GLMM)*

22 GLMMs were used to detect SNPs correlated with environmental variables as
23 previously described [31]. A GLMM model with logit link and a binomial error distribution
24 was estimated between each SNP and the environmental variable. The R lme4 package was
25 used to calculate the models [32]. Using the allelic frequencies of the alleles at a defined

1 genomic position, the number of occurrences of each allele was calculated for a total count of
2 30, and these occurrence data were used as input for model determination. This strategy enabled
3 us to perform such calculations based on an equal total count of alleles at each location at a
4 genomic position. To evaluate model performances, the likelihood ratio (LR) was calculated
5 using the ANOVA function from the R car package with the model and null model as input. A
6 Bonferroni correction [33] was applied to the LR p-values as previously described [31]. Wald1
7 and Wald2 coefficients were also used to evaluate model performances and were recovered
8 from the GLMM model data produced [34].

9

10 *2.5.3 Latent factor mixed models (LFMM)*

11 The LFMM 1.5 program enables testing for the association of genomic loci with
12 environmental variables while taking into account the population structure through unobserved
13 latent factors ([35], <http://membres-timc.imag.fr/Olivier.Francois/lfmm/index.htm>). The
14 number of latent factors is fixed using the K parameter. To define the optimal value, we tested
15 values from 2 to 6 with a reduced dataset containing a total of 120,000 SNPs, 15000 per
16 chromosome. Genotyping data were provided as allelic frequencies per population. LFMM
17 implement an MCMC algorithm, the number of iterations in the Gibbs sampling algorithm was
18 set to 30,000, and the number of burning iterations to 15,000. For the complete dataset, we
19 performed 6 independent runs, and the Z-scores were combined using the Fisher-Stouffer
20 method. The p-values were adjusted as described by LFMM manual ([http://membres-](http://membres-timc.imag.fr/Olivier.Francois/lfmm/index.htm)
21 [timc.imag.fr/Olivier.Francois/lfmm/index.htm](http://membres-timc.imag.fr/Olivier.Francois/lfmm/index.htm)).

22

23

24

25

2. Results and discussion

2.1 Genomic sequencing and population structure analysis

In order to acquire genomic variation data from populations of *Arabis alpina* growing in contrasting environments, we sampled leaves from individual plants along an elevation gradient in a range of 860 to 2,781 m (Table 1).

Extracted DNA was submitted to pool sequencing to get the allelic frequencies at a large number of loci across the genome, each pool containing an equivalent amount of 10 independent individuals (Supplementary data, Table S1 for sequencing statistics). A sampled population contains between 20 and 30 individual plants depending on the sampling sites, with only plants that are at least 1 meter apart being sampled since *Arabis alpina* tends to propagate at short distances through sucker production. Ten sampled plants represent therefore a significant part of each studied population, suitable for allelic frequency estimation. SNP mapping using the 8 *Arabis alpina* chromosomes that have been sequenced as a reference [20] produced a dataset of 2,575,324 SNPs for which data are available in all sampled populations. This number of SNPs corresponded to a density of one SNP at every 130 bp and well covered the whole sequenced genome of *Arabis alpina*. A phylogenetic tree was built from a subset of 8,000 SNPs to detect the genetic relationship between the sampled populations (Fig. 1). A clear genetic structure was detected, and 3 groups of populations were identified. Two groups correspond to a set of samples from the Chartreuse and Vercors area, that are pre-alpine massifs of the French Alps. A third group corresponds to samples collected in the French Alps. Interestingly, samples in the lower and higher range of elevation of the altitude gradient were collected in all 3 areas, which may help to prevent the detection of false positives when looking for loci associated with climatic variables by limiting population effects.

1 2.2 Combined genome scan

2

3 To detect associations between loci as SNPs and the climatic environment, we prepared
4 a synthetic variable summarizing the climatic conditions in which these plant populations grow.
5 For this purpose, a PCA was performed using 3 available variables (annual range of
6 temperatures, annual precipitations, and snow coverage). The first axis of the PCA explained
7 97% of the variance of the data and was chosen as the climatic variable for the *Arabis alpina*
8 sampling sites. Two independent analyses were then performed with statistical tools that enable
9 us to detect associations between loci and the environmental variables.

10

11 We first used GLMM, which we have previously been shown to be effective in detecting
12 associations between minimum temperature and SNPs in *Arabidopsis* [31]. Indeed, the analysis
13 of 80 whole-genome data from 80 *Arabidopsis* ecotypes allowed identifying several significant
14 loci, most of them documented as being involved in stress or cold response in *Arabidopsis*
15 *thaliana*. We generated a model for each *Arabis alpina* SNP of the dataset, using the 3 groups
16 identified in the phylogenetic tree as population structure data to take into account that
17 parameter and reduce the detection of false positives. The presence of low and high elevation
18 sites in each group may also contribute to reducing false positives through the population effect.

19 SNPs that were significantly associated with the synthetic climatic variable used in this
20 study were selected using LR and a Wald test according to the procedure previously described
21 [31]. At a threshold of 0.01 and after a Bonferroni correction, 533,699 SNPs were considered
22 as positive. The lowest LR p-value in this analysis was $1.497 \cdot 10^{-83}$, well below the threshold
23 used.

24

1 In a second approach, we used the LFMM program to test for the association of SNP
2 allelic frequency with the climatic variable. Six independent analyses were performed and
3 combined as suggested by the program developers [35]. In this case, the software itself analyzed
4 the population structure and included this parameter in the analysis. The K latent factor was
5 fixed at 3 as the optimal value, which was in agreement with the 3 groups of populations
6 identified in the phylogenetic analysis (Fig. 1). A 0.01 threshold led to the identification of
7 141,392 positive SNPs, with a lowest adjusted p-value of $2.7 \cdot 10^{-9}$.

8
9 In order to extract from these two analyses the most significantly associated SNPs, and
10 considering the large number of putative positive SNPs detected at the thresholds used , we
11 selected the 500 SNPs corresponding to the lowest p-values from each method. We then
12 combined the two sets of results. The SNPs detected by both methods were selected.
13 Subsequently, those SNPs were retained that correspond to genomic regions for which at least
14 two positive SNPs were detected within a 5kb window. Such a sliding-window analysis was
15 successfully used to detect SNPs associated with minimum temperature in *Arabidopsis thaliana*
16 [31].

17

18 *2.3 Analysis of detected genomic regions*

19 *2.4 A major genomic region located on chromosome 5*

20 Among the 19 regions detected (Supplementary data, Table S2), the most significant
21 was located on *Arabis alpina* chromosome 5. All the 25 first LFMM lowest p-values are located
22 within a 25kb genomic region of this chromosome, from position 11,521,406 to 11,546,695.
23 Furthermore, this part of chromosome 5 was also the most significant one with the GLMM
24 approach, with 13 hits within the 21st lowest p-values generated by this method. Therefore, a

1 strong positive signal was detected with both methods. This genomic region harbors 4 annotated
2 genes (Fig. 2).

3 Aa_G315180, described as the ortholog of the *Arabidopsis thaliana* At3g49490 gene,
4 encodes an uncharacterized protein.

5 Aa_G315190 and Aa_G315200 are annotated as pyruvate decarboxylase genes (PDC).
6 In *Arabidopsis*, a family of 3 *PDC* genes has been described [36]. They encode proteins of 592
7 to 607 amino acids long, which is also the average size for this enzyme in other plants. The
8 putative *Arabis PDC1* and 2 genes, however, encode much shorter proteins of respectively 111
9 and 306 aa. When compared to the structure of the *Arabidopsis PDC* genes, the *PDC1* encoding
10 region is interrupted by a stop codon in the 4th exon. For *PDC2*, the whole gene structure is
11 detected in this part of the *Arabis* chromosome 5; however, a stop codon leads to a truncated
12 protein. We have previously assembled *Arabis alpina* genomic sequences, and the contig
13 7180002427937 of this assembly corresponds to this part of chromosome 5. These data confirm
14 the sequence documented in the *Arabis alpina* genome [18], and *Arabis PDC1* and 2 are
15 therefore very likely to be pseudogenes.

16

17 Aa_G315210 is annotated as the ortholog of the *Arabidopsis thaliana* gene At3g49530,
18 which encodes a transcription factor named NAC062. The structure of the two genes is similar,
19 with 6 exons and the corresponding proteins of 468 and 457 amino acid sharing 77% identity
20 at the amino acid level. Among the SNPs detected in the Aa_G315210 gene, most of them
21 correspond to synonymous mutations or are located in non-coding regions, and one of them is
22 a non-synonymous mutation. At position 182 of the protein, the two alleles lead to a shift from
23 a proline to an arginine (Fig. 3). This SNP represents therefore a non-conservative substitution
24 in the *Arabis* NAC062 protein and its allelic profile along the elevation gradient was detected
25 as highly correlated to the synthetic variable used. Alignment of *Arabidopsis* and *Arabis*

1 NAC062 protein sequences revealed that position 182 of the *Arabis* protein corresponds to an
2 amino acid insertion compared to the counterpart in *Arabidopsis*. This is also true when
3 *AaNAC062* is compared with its ortholog from *Capsella rubra* or *Arabidopsis lyrata*, other
4 plants from the Brassicaceae family. It was therefore not possible to investigate whether similar
5 mutations at this position occur in the different *Arabidopsis* ecotypes for which genomic data
6 are available. Fig. 4 shows an example of an allelic profile that we detected as being highly
7 correlated with the environmental variable within the Aa_G315210 gene. A clear shift in allelic
8 frequency is detected along the gradient, with distinct alleles present at the bottom and the top
9 part of the climatic gradient.

10

11 The Aa_G315210 gene encodes a protein of the MTF family (membrane-bound
12 transcription factors) annotated as the ortholog of the NAC062/NTL6 protein of *Arabidopsis*
13 *thaliana*. MTFs are membrane-anchored proteins that respond to stimulation by migrating to
14 the nucleus, where they regulate the expression of their target genes [37]. Such transcription
15 factors are involved in the stress response of different organisms [38], allowing a rapid
16 transcriptional regulation under changing environmental conditions [39, 40]. In *Arabidopsis*,
17 MTFs have been shown to be involved in the response to different stress factors such as high
18 salt and cold [41-43] and developmental signals [44, 45].

19 The AtNAC062 protein is activated in response to cold. A primary event upon exposure
20 to low temperatures is a change in membrane fluidity that induces the proteolytic processing of
21 NAC062 and its release from the plasma membrane [46, 47]. NAC062 activates the expression
22 of a set of PR proteins (pathogenesis-related proteins), linking the cold stress and the pathogen
23 responses [41].

24

1 In winter, snow acts as a protection against freezing temperatures but maintains also a
2 high humidity favorable for psychrophilic pathogens to attack dormant plants [48, 49]. This has
3 been shown for overwintering plants that need to keep vegetative tissues under the snow. Long-
4 term snow cover leads to the depletion of plant carbohydrate reserves, so plants become more
5 susceptible to pathogen attack [50]. In these overwintering species, PR proteins protect against
6 pathogen attacks and also provide freezing tolerance [51]. The perennial plant *Arabis alpina*
7 also keeps vegetative tissue under the snow cover to allow for rapid flowering and seed
8 production after the snow has melted in the spring at high elevation [52]. Along the studied
9 climatic gradient, periods of growth without snow varies in a large range from 9 to 42 weeks at
10 high elevations and may lead to specific adaptation to face this constraint.

11

12 2.5 Other genomic regions associated to the synthetic climatic variable

13 Among the other regions detected as being significantly associated with the synthetic
14 climatic variable used in this study, some contain genes linked to the stress response (Table 2).
15 The Aa_G89730 gene encodes a protein of the pfkb-like carbohydrate kinase family. Such
16 proteins have been shown to be involved in the drought response in *Phaseolus vulgaris* [53].
17 On chromosome 8, an acyl-activating enzyme gene is located within the detected region. This
18 family of enzymes perform carboxylic acid activation and is required in different metabolic
19 pathways [54]. In *Arabidopsis*, Rizhsky *et al.* [55] have shown, through transcriptome analysis,
20 that the expression of genes of this family is modulated by a combination of abiotic stress
21 factors. They exposed plants simultaneously to drought and heat. Such situations of exposure
22 to multiple stress factors are very likely to occur in mountain areas. Gene Aa_G18650 encodes
23 a heat shock protein.

24 An ortholog of the *Arabidopsis thaliana* *AtABI3* gene was detected on
25 chromosome 3 (abscisic acid-insensitive 3). ABI3 protein has been initially characterized as a

1 transcription factor involved in plant seed maturation. But recently some other functions have
2 been reported, such as a role in plant development and flowering time [56]. Interestingly,
3 another significant region contains an ortholog of the *AtPRMT10* gene on *Arabis alpina*
4 chromosome 8. It encodes a protein called arginine methyltransferase, which catalyzes the
5 asymmetric demethylation of arginine 3 of histone 4. This enzyme is involved in the regulation
6 of flowering time [57]. The *Arabidopsis* PRMT10 mutant has a late flowering time phenotype.
7 In *Arabis alpina*, flowering is linked to cold since a prolonged exposure to low temperatures is
8 required to initiate the formation of flower buds [58].

9
10 Members of other gene families were also detected. An aluminum-activated malate
11 transporter gene, encoding a protein of a family that has first been characterized in wheat as
12 being involved in aluminum resistance, but these proteins perform many other functions as well
13 [59]. A pectate lyase encoding enzymes important in various plant developmental processes
14 was also detected [60]. Finally, some genes for which no functional data are available so far,
15 such as hypothetical proteins or the ortholog of an *Arabidopsis* kinase-like protein (TMKL1)
16 for which no cellular function has been described yet [61], were also found. Some regions do
17 not contain annotated genes according to the actual knowledge of the *Arabis alpina* genome,
18 but we cannot rule out that they may play a role in *Arabis alpina* adaptation.

19

20 **3. Discussion**

21 Genomic studies have enabled investigating local adaptation at the molecular level by
22 looking for genome-environment associations in plant species [3]. Major contributions have
23 been made using the model plant *Arabidopsis thaliana*, for which abundant resources are
24 available [62-65]. In addition to sequenced genomes, functional data and characterized
25 ecotypes from different area allow detailed investigations. The recent evolution in sequencing

1 techniques enables to generate high-level genomic resources in non-model organisms, opening
2 the possibility of studying local adaptation in other species using high-density data of genetic
3 variations. Although *Arabidopsis thaliana* has been useful in providing some knowledge of
4 local adaptation in plants, other models may enable addressing specific questions that cannot
5 be investigated in *Arabidopsis*. *Arabis alpina* has become a complementary species to the
6 annual plant *Arabidopsis* to study flowering time regulation in a perennial plant [22]. Since the
7 two plants are closely related in the Brassicaceae family, *Arabis alpina* studies benefit from
8 functional gene data available from *Arabidopsis thaliana* resources. In this work, we
9 investigated this species as a suitable model for studying local adaptation of alpine plants.
10 Temperature is a well-known key factor for plant growth, and association studies with climatic
11 variables in *Arabidopsis* have confirmed that it is an essential parameter [62, 65]. The
12 increasing duration of the growing season, which in part is correlated with local temperature,
13 is also an important factor for *Arabidopsis* [64]. The synthetic variable that we used in our study
14 integrates these two factors with the annual precipitations.

15 If these variables are common drivers in local adaptation of plants, the major loci that
16 we identified appear to be linked to the life cycle of a perennial alpine plant growing in
17 contrasting environments along an elevation gradient. For positive SNPs associated with the
18 synthetic environmental variable, we detected a shift in alleles along the elevation gradient.
19 These results suggest that environmental pressures indeed shape genetic diversity along the
20 gradient, leading to local adaptations. Alpine plants are mainly perennials and adapt to the
21 specific climatic conditions where they grow by initiating flowering in response to cold through
22 vernalization [58, 66]. In this way, flowering is restricted to spring at the beginning of the short
23 growing season so that seeds and next plantlets are produced with optimal chances of survival.
24 Apart from genes linked to stress responses, some major candidate genes that we identified are
25 potentially involved in adaptations to variations in snow cover and therefore growing period;

1 these variations may be involved in plant protection under long-term snow cover. Flowering
2 time is also crucial and is linked to low temperatures in *Arabis* [58], and candidate loci
3 associated with flowering regulation were detected. Some candidate genes are therefore
4 strongly related to key drivers of plant adaptation in mountain areas. Selective pressure
5 associated with local climate acting through temperature and growing season varies along the
6 studied gradient, with an increase in selective pressure at high elevation.

7 *Arabis alpina*, therefore, represents a suitable model for studying genomic adaptation
8 in alpine perennial plants, for exploring adaptation to climate or other environmental factors,
9 and for investigating the functions of candidate gene variants.

10

11

12

13

14

15

16 **Acknowledgments**

17 We are grateful to Marie Aline for help in *Arabis alpina* sampling. Many thanks to
18 Mathieu Gautier and Michael Blum for discussions about this work.

19

20

21

22

23

24

25

1 **References**

2

- 3 1. J. Hereford, A quantitative survey of local adaptation and fitness trade-offs, *Am. Nat.*
4 173 (2009) 579-588.
- 5 2. O. Savolainen, M. Lascoux, J. Merilä, Ecological genomics of local adaptation, *Nat.*
6 *Rev. Genet.* 14 (2013) 807-820.
- 7 3. C. Weinig, B.E. Ewers, S.M. Welch, Ecological genomics and process modeling of
8 local adaptation to climate, *Curr. Opin. Plant. Biol.* 18 (2014) 66-72.
- 9 4. S.D. Schoville, et al., Adaptive Genetic Variation on the Landscape: Methods and
10 Cases, *Annu. Rev. Ecol. Evol. Syst.* 43 (2012) 23-43.
- 11 5. C.W. Ahrens, et al., The search for loci under selection: trends, biases and progress,
12 *Mol. Ecol.* 27 (2018) 1342-1356.
- 13 6. S. Hoban, et al., Finding the Genomic Basis of Local Adaptation: Pitfalls, Practical
14 Solutions, and Future Directions, *Am. Nat.* 188 (2016) 379-397.
- 15 7. S. Goodwin, J.D. McPherson, W.R. McCombie, Coming of age: ten years of next-
16 generation sequencing technologies, *Nat. Rev. Genet.* 17 (2016) 333-351.
- 17 8. R.J. Elshire, et al., A robust, simple genotyping-by-sequencing (GBS) approach for high
18 diversity species, *PLoS One* 6 (2011) e19379.
- 19 9. N.A. Baird, et al., Rapid SNP discovery and genetic mapping using sequenced RAD
20 markers, *PLoS One* 3 (2008) e3376.
- 21 10. A.P. Fuentes-Pardo, D.E. Ruzzante, Whole-genome sequencing approaches for
22 conservation biology: Advantages, limitations and practical recommendations, *Mol.*
23 *Ecol.* 26 (2017) 5369-5406.

- 1 11. C. Rellstab, et al., Validation of SNP allele frequencies determined by pooled next-
2 generation sequencing in natural populations of a non-model plant species, PLoS One
3 8 (2013) e80422.
- 4 12. Y. Zhu, et al., Empirical validation of pooled whole genome population re-sequencing
5 in *Drosophila melanogaster*, PLoS One 7 (2012) e41901.
- 6 13. M. Ozerov, et al., Cost-effective genome-wide estimation of allele frequencies from
7 pooled DNA in Atlantic salmon (*Salmo salar* L.), BMC Genomics 14 (2013) 12.
- 8 14. M.A. Fustier, et al., Signatures of local adaptation in lowland and highland teosintes
9 from whole-genome sequencing of pooled samples, Mol. Ecol. 26 (2017) 2738-2756.
- 10 15. K.N. Stölting, et al., Genome-wide patterns of differentiation and spatially varying
11 selection between postglacial recolonization lineages of *Populus alba* (Salicaceae), a
12 widespread forest tree, New Phytol. 207 (2015) 723-34.
- 13 16. C. Körner, The use of 'altitude' in ecological research, Trends Ecol. Evol. 22 (2007)
14 569-574.
- 15 17. C. Körner, *Alpine Plant Life: Functional Plant Ecology of High Mountain Ecosystems*,
16 2003, Springer, Heidelberg.
- 17 18. S. Lobréaux, S. Manel, C. Melodelima, Development of an *Arabis alpina* genomic
18 contig sequence data set and application to single nucleotide polymorphisms discovery.
19 Mol. Ecol. Resour. 14 (2014) 411-418.
- 20 19. E.M. Willing, et al., Genome expansion of *Arabis alpina* linked with retrotransposition
21 and reduced symmetric DNA methylation, Nat. Plants 1 (2015) 14023.
- 22 20. W.B. Jiao, et al., Improving and correcting the contiguity of long-read genome
23 assemblies of three plant species using optical mapping and chromosome conformation
24 capture data, Genome Res. 27 (2017) 778-786.

- 1 21. S. Bergonzi, et al., Mechanisms of age-dependent response to winter temperature in
2 perennial flowering of *Arabis alpina*, *Science* 340 (2013). 1094-1097.
- 3 22. R. Wang, et al., PEP1 regulates perennial flowering in *Arabis alpina*, *Nature* 459 (2009)
4 423-427.
- 5 23. B.N. Poncet, et al., Tracking genes of ecological relevance using a genome scan in two
6 independent regional population samples of *Arabis alpina*, *Mol. Ecol.* 19 (2010). 2896-
7 2907.
- 8 24. D. Zulliger, E. Schnyder, F. Gugerli, Are adaptive loci transferable across genomes of
9 related species? Outlier and environmental association analyses in Alpine Brassicaceae
10 species, *Mol. Ecol.* 22 (2013) 1626-1639.
- 11 25. H. Li, R. Durbin, Fast and accurate short read alignment with Burrows-Wheeler
12 transform, *Bioinformatics* 25 (2009) 1754-1760.
- 13 26. C. Melodelima, S. Lobréaux, Complete *Arabis alpina* chloroplast genome sequence and
14 insight into its polymorphism, *Meta Gene* 1 (2013) 65-75.
- 15 27. H. Li, et al., The Sequence Alignment/Map format and SAMtools, *Bioinformatics* 25
16 (2009) 2078-2079.
- 17 28. N. Takezaki, M. Nei, K. Tamura, POPTREEW: web version of POPTREE for
18 constructing population trees from allele frequency data and computing some other
19 quantities, *Mol. Biol. Evol.* 31 (2014) 1622-1624.
- 20 29. A.A. Popescu, K.T. Huber, E. Paradis, ape 3.0: New tools for distance-based
21 phylogenetics and evolutionary analysis in R, *Bioinformatics* 28 (2012) 1536-1537.
- 22 30. D. Chessel, A.B. Dufour, J. Thioulouse, The ade4 Package – I: One-Table Methods, *R*
23 *News* 4 (2004) 5-10.
- 24 31. S. Lobréaux, C. Melodelima, Detection of genomic loci associated with environmental
25 variables using generalized linear mixed models, *Genomics* 105 (2015) 69-75.

- 1 32. D. Bates, M. Mächler, B. Bolker, S. Walker, Fitting Linear Mixed-Effects Models Using
2 lme4, *J. Stat. Soft.* 67 (2015) 1-48.
- 3 33. J.P. Shaffer, Multiple hypothesis-testing, *Ann. Rev. of Psychol.* 46 (1995) 561-584.
- 4 34. S. Joost, A spatial analysis method (SAM) to detect candidate loci for selection: towards
5 a landscape genomics approach to adaptation, *Mol. Ecol.* 16 (2007) 3955-3969.
- 6 35. E. Frichot, E., et al., Testing for associations between loci and environmental gradients
7 using latent factor mixed models, *Mol. Biol. Evol.* 30 (2013) 1687-1699.
- 8 36. O. Kürsteiner, I. Dupuis, C. Kuhlemeier, The pyruvate decarboxylase1 gene of
9 Arabidopsis is required during anoxia but not other environmental stresses, *Plant*
10 *Physiol.* 132 (2003) 968-978.
- 11 37. S.Y. Kim, et al., Exploring membrane-associated NAC transcription factors in
12 Arabidopsis: implications for membrane biology in genome regulation, *Nucleic Acids*
13 *Res* 35 (2007) 203-213.
- 14 38. P.J. Seo, Recent advances in plant membrane-bound transcription factor research:
15 emphasis on intracellular movement, *J. Integr. Plant Biol.* 56 (2014) 334-342.
- 16 39. Y.N. Chen, E. Slabaugh, F. Brandizzi, Membrane-tethered transcription factors in
17 Arabidopsis thaliana: novel regulators in stress response and development, *Curr. Opin.*
18 *Plant Biol.* 11 (2008) 695-701.
- 19 40. P.J. Seo, S.G. Kim, C.M. Park, Membrane-bound transcription factors in plants, *Trends*
20 *Plant Sci.* 13 (2008) 550-556.
- 21 41. P.J. Seo, et al., Cold activation of a plasma membrane-tethered NAC transcription factor
22 induces a pathogen resistance response in Arabidopsis, *Plant J.* 61 (2010) 661-671.
- 23 42. H.K. Yoon, et al., Regulation of leaf senescence by NTL9-mediated osmotic stress
24 signaling in Arabidopsis, *Mol. Cells* 25 (2008) 438-445.

- 1 43. M.J. Kim, et al., Controlled nuclear import of the transcription factor NTL6 reveals a
2 cytoplasmic role of SnRK2.8 in the drought-stress response, *Biochem J.* 448 (2012)
3 353-363.
- 4 44. Y.S. Kim, et al., A membrane-bound NAC transcription factor regulates cell division in
5 *Arabidopsis*, *Plant Cell* 18 (2006) 3132-3144.
- 6 45. S.G. Kim, et al., A membrane-bound NAC transcription factor NTL8 regulates
7 gibberellic acid-mediated salt signaling in *Arabidopsis* seed germination, *Plant J.* 55
8 (2008) 77-88.
- 9 46. B.L. Orvar, et al., Early steps in cold sensing by plant cells: the role of actin
10 cytoskeleton and membrane fluidity, *Plant J.* 23 (2000) 785-794.
- 11 47. P.J. Seo, et al., Proteolytic processing of an *Arabidopsis* membrane-bound NAC
12 transcription factor is triggered by cold-induced changes in membrane fluidity,
13 *Biochem J.* 427 (2010) 359-367.
- 14 48. M. Griffith, M.W. Yaish, Antifreeze proteins in overwintering plants: a tale of two
15 activities, *Trends Plant Sci.* 9 (2004) 399-405.
- 16 49. C.S. Snider, et al., Role of ice nucleation and antifreeze activities in pathogenesis and
17 growth of snow molds, *Phytopathology* 90 (2000) 354-361.
- 18 50. T. Nakajima, J. Abel, Development of resistance to *Microdochium nivale* in winter
19 wheat during autumn and decline of the resistance under snow, *Can J. Bot.* 72 (1994)
20 1211-1215.
- 21 51. W.C. Hon, et al., Antifreeze proteins in winter rye are similar to pathogenesis-related
22 proteins, *Plant Physiol.* 109 (1995) 879-889.
- 23 52. P. Toräng, et al., Large-scale adaptive differentiation in the alpine perennial herb *Arabis*
24 *alpina*, *New Phytol.* 206 (2015) 459-470.

- 1 53. K. Tatjana, et al., Identification of genes involved in the response of leaves of
2 *Phaseolus vulgaris* to drought stress, *Molecular Breeding* 21 (2008) 159:172.
- 3 54. J.M. Shockey, M.S. Fulda, J. Browse, *Arabidopsis* contains a large superfamily of acyl-
4 activating enzymes. Phylogenetic and biochemical analysis reveals a new class of acyl-
5 coenzyme a synthetases, *Plant Physiol.* 132 (2003) 1065-1076.
- 6 55. L. Rizhsky, et al., When defense pathways collide. The response of *Arabidopsis* to a
7 combination of drought and heat stress, *Plant Physiol.* 134 (2004) 1683-1396.
- 8 56. A. Rohde, S. Kurup, and M. Holdsworth, ABI3 emerges from the seed, *Trends Plant.*
9 *Sci.* 5 (2000) 418-419.
- 10 57. L. Niu, L., et al. Regulation of flowering time by the protein arginine methyltransferase
11 AtPRMT10, *EMBO Rep.* 8 (2007) 1190-1195.
- 12 58. A. Lazaro, E. Obeng-Hinneh, M.C. Albani, Extended Vernalization Regulates
13 Inflorescence Fate in *Arabis alpina* by stably silencing Perpetual Flowering1, *Plant*
14 *Physiol.* 176 (2018) 2819-2833.
- 15 59. A.J. Palmer, A. Baker, S.P. Muench, The varied functions of aluminium-activated
16 malate transporters-much more than aluminium resistance, *Biochem. Soc. Trans.* 44
17 (2016) 856-862.
- 18 60. S.G. Palusa, et al., Organ-specific, developmental, hormonal and stress regulation of
19 expression of putative pectate lyase genes in *Arabidopsis*, *New Phytol.* 174 (2007) 537-
20 550.
- 21 61. C. Valon, C., et al., Characterization of an *Arabidopsis thaliana* gene (TMKL1)
22 encoding a putative transmembrane protein with an unusual kinase-like domain, *Plant*
23 *Mol. Biol.* 23 (1993) 415-421.
- 24 62. A. Fournier-Level, et al., A map of local adaptation in *Arabidopsis thaliana*, *Science*
25 334 (2011) 86-89.

1 63. A.M. Hancock, et al., Adaptation to climate across the *Arabidopsis thaliana* genome,
2 Science 334 (2011) 83-86.

3 64. J.R. Lasky, et al., Characterizing genomic variation of *Arabidopsis thaliana*: the roles
4 of geography and climate, *Mol. Ecol.* 21 (2012) 5512-5529.

5 65. Y. Zhen, M.C. Ungerer, Clinal variation in freezing tolerance among natural accessions
6 of *Arabidopsis thaliana*, *New Phytol.* 177 (2008) 419-427.

7 66. A. Aydelotte, P. Diggle, Analysis of developmental preformation in the alpine herb
8 *Caltha leptosepala* (Ranunculaceae), *Am. J. Bot.* 84 (1997) 1646.

9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

1 **Figure legends**

2

3 **Fig. 1.** Phylogenetic tree of the ten *Arabis alpina* populations studied based on SNP
4 data. A phylogenetic tree was built using the neighbor-joining method from a set of 8,000 SNPs
5 by randomly selecting 1,000 SNPs from each *Arabis alpina* chromosome. The three major
6 groups were detected in the tree and are indicated as A, B, and C.

7

8 **Fig. 2.** The region containing *AaNAC062* gene is a major locus associated with the
9 environmental climatic variable. A) Schematic representation of the chromosome-5 genomic
10 region from position 11,521,406 to 11,546,695 bp. The different genes are indicated by an
11 arrow from ATG to stop codons, with names according to the *Arabis alpina* genome
12 annotations. Black boxes correspond to exons. Intergene sequences are represented as dashed
13 lines, and their length is indicated.

14

15 **Fig. 3.** Partial alignment of NAC062 protein from *Arabis alpina* (AaG315210.t1),
16 *Arabidopsis thaliana* NAC062 (AT3G49530.1), *Arabidopsis lyrata subsp lyrata*
17 (XP_002877683.1) and *Capsella rubella* (XP_006291082.1) from amino acids 171 to 189. The
18 genetic variation detected between *Arabis alpina* populations in the climatic gradient, and
19 responsible for a non-synonymous mutation, is indicated by a star at position 181.

20

21 **Fig. 4.** Plot of the allelic frequency profile for a sample SNP from *AaNAC062* gene
22 against the synthetic environmental variable used in the study. The upper part of the gradient at
23 high elevation corresponds to positive values of the variable.

24

25

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16

Table 1

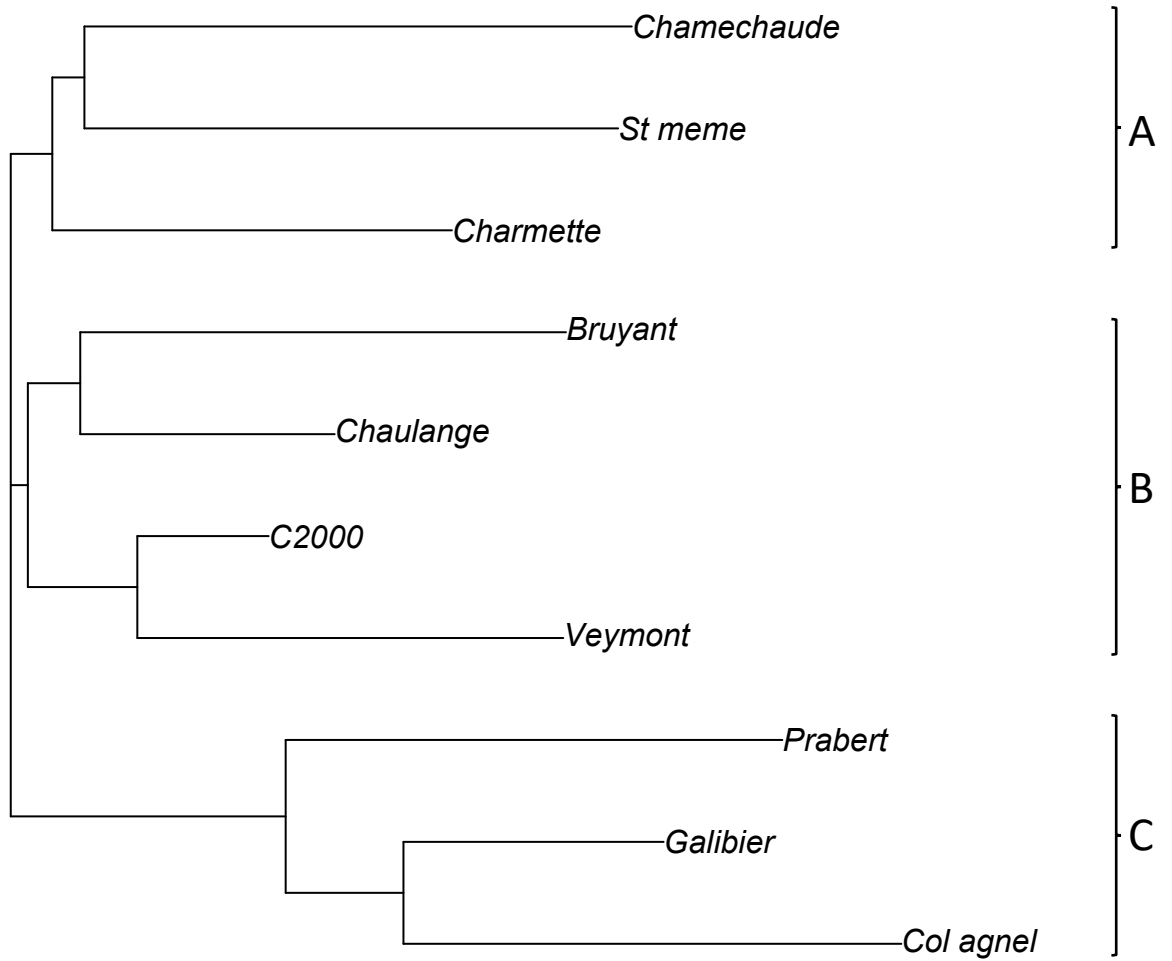
Sampling site data and synthetic climatic variable corresponding to axis 1 values from PCA of environmental variables

Name	Latitude	Longitude	Elevation (m)	Synthetic variable
St meme	45.39750	5.89125	860	-2.4427
Bruyant	45.15069	5.61153	980	-1.4517
Charmette	45.32818	5.73670	1364	-0.4670
Prabert	45.21890	5.99395	1380	-1.4145
Chaulange	45.07060	5.59411	1545	-1.0925
Chamechaude	45.28571	5.78504	1914	0.6757
Veymont	44.88069	5.52211	1950	0.4184
C2000	45.01968	5.57005	1951	0.1420
Galibier	45.06058	6.40485	2536	2.9396
Col agnel	44.68568	6.98280	2781	2.6927

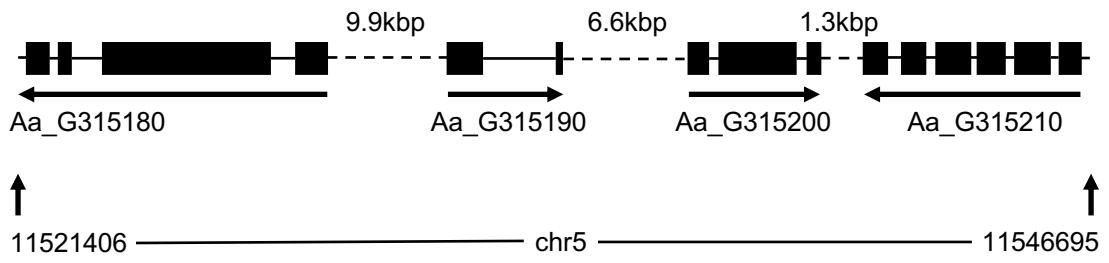
1 **Table 2**

2 *Arabis alpina* candidate genes within genomic regions significantly associated with the synthetic environmental
 3 variable. Potential candidate genes within these regions are indicated together with their *Arabidopsis thaliana*
 4 ortholog when available.

5	Chr	localization	Aa_ID	Function	Arabidopsis ortholog
6	chr1	1132833-1137852	Aa_G6950	uncharacterized protein	AT1G02990
7	chr1	4989967-4992419	Aa_G368310	aluminum activated malate transporter	AT1G08440
8	chr2	35950629-35952085	Aa_G797250	hxxxd-type acyl-transferase-like protein	AT1G78990
9	chr3	1816190-1819332	Aa_G185650	tcp-1 cpn60 chaperonin family protein	AT3G03960
10	chr3	17087273-17089595	Aa_G284660	uncharacterized protein	AT3G24750
11	chr3	17108584-17111085	Aa_G284690	pectate lyase family protein	AT3G24670
12	chr3	17111534-17113646	Aa_G284700	kinase-like protein tmkl1-like	AT3G24660
13	chr3	17136875-17137113	Aa_G385880	leucine-rich repeat-containing protein	-
14	chr3	17140832-17143771	Aa_G385870	abscisic acid-insensitive protein 3	AT3G24650
15	chr5	11544865-11546695	Aa_G315210	NAC062 transcription factor	AT3G49530
16	chr5	12159593-12705643	Aa_G53000	Kinesin family-like protein	-
17	chr6	2380152-2382801	Aa_G89730	pfkb-like carbohydrate kinase protein	AT5G43910
18	chr6	30161491-30162716	Aa_G245490	vesicle transport v-snare	AT5G39510
19	chr8	35675667-35677688	Aa_G54570	protein arginine n-methyltransferase	AT1G04870
20	chr8	45693680-45697947	Aa_G97990	acyl-activating enzyme 17	AT5G23050
21					



0.1



*

Aa_NAC062	171	EEDSKSDEAEPEPAGSSPT	189
		:	
At_NAC062	171	EEDSKSDEVE-EPAVSSPT	188

