



GIO – Genes In Order – A suite of scripts for identification and mapping of SNP markers in GBS data

Journal:	<i>The Plant Genome</i>
Manuscript ID	TPG-2015-10-0094
Manuscript Type:	Science Notes
Date Submitted by the Author:	01-Oct-2015
Complete List of Authors:	Skinner, Daniel; USDA-ARS, Crop and Soil Sciences Krishnan, Vandhana; Washington State University, Crop and Soil Sciences See, Deven; USDA-ARS and Washington State University, Plant Pathology
Keywords:	Genome sequencing, Genotyping by sequencing, Next generation sequencing

SCHOLARONE™
Manuscripts

Only

1
2
3 **GIO – Genes In Order – A suite of scripts for identification and mapping of SNP markers**
4 **in GBS data**
5
6

7 **Daniel Z. Skinner***, Vandhana Krishnan, and Deven R. See
8

9
10 D.Z. Skinner and D.R. See, U.S. Department of Agriculture, Agricultural Research Service, 209
11 Johnson Hall, Washington State University, Pullman, WA 99164, U.S.A. V. Krishnan and D.Z.
12 Skinner, Department of Crop and Soil Sciences, 115 Johnson Hall, Washington State University,
13 Pullman, WA 99164, U.S.A. D.R. See, Department of Plant Pathology, 345 Johnson Hall,
14 Washington State University, Pullman, WA 99164, U.S.A. *Corresponding author
15
16
17
18 (dan.skinner@ars.usda.gov).
19

20
21 Abbreviations: NGS, Next Generation Sequencing; SNP, Single nucleotide Polymorphism; GIO,
22 Genes In Order.
23

24
25 **Abstract**
26

27
28 Next generation sequencing (NGS) generates billions of short DNA sequencing “tags” that can
29 be used to identify genomic regions associated with trait determination through inheritance or
30 genome-wide association studies. First, it is necessary to identify pairs of tags that differ by
31 nucleotide base(s), tabulate occurrence among sample lines, and infer genomic locations in
32 mapping populations. GIO provides a simple-to-use tool to carry out these functions through a
33 simple graphical user interface. The output from an association panel study consists of the
34 putative chromosome of origin and the presence/absence of each tag in the sample lines, and an
35 indication of whether the tag was homologous to known expressed sequences. A separate table
36 of tags that differed from a single other tag in the dataset by the user-specified number of
37 SNP(s), and thus may represent alleles, is produced including the genotypes of the sample lines.
38 The output from a biparental study also includes progeny genotypes, separated into low-copy and
39 high copy-number tags. Two draft genetic maps are produced, one based on a user-specified
40 number of genetic distance clusters, and one based on homology to previously-mapped
41 sequences. Thousands of dominant tags typically are identified in biparental studies and GIO
42 provides their putative chromosome of origin, and presence/absence in the progeny. A table of
43 pairs of tags (both co-dominant and dominant) that cosegregate with at least one other tag in at
44 least 65% of the progeny is also provided to facilitate mapping of the dominant tags. GIO is
45 available from the first author.
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 **Keywords:** Next generation sequencing, Computational genomics; Genotyping; mapping;
4 Recombination; SNPs
5

6
7 Next generation sequencing has provided a means to quickly generate millions of short DNA
8 sequence reads from study subjects with no required prior knowledge of the genome. The
9 potential benefit of the application of these “tags” in population studies to link genetic variation
10 to phenotypic variation is enormous. Variant forms of stretches of DNA that differ by a single
11 nucleotide base (single nucleotide polymorphism, or “SNP”) can be found throughout the
12 genome, theoretically resulting in DNA sequence variation occurring within a short distance of
13 all regulatory and coding regions (Sachidanandam et al., 2001). These SNPs can be used as
14 biological markers, helpful to locating genes that are associated with phenotypic traits of interest
15 (Lander and Botstein, 1989). The tags generated by high-throughput sequencing are a highly
16 efficient means of identifying genomic regions containing SNPs.
17

18
19 In order to identify useful markers, the DNA sequence of each of the DNA tags generated by
20 high-throughput sequencing must be compared to the sequences of the millions of other tags
21 produced. Duplicate tags must be identified and removed, and useful variant forms identified in
22 the appropriate proportion of individuals in the study. By comparing the tags from one parent to
23 tags from the other parent in a biparental population study, tags representing co-dominant alleles
24 may be discovered, as well as tags that do not have a corresponding alternate allele represented.
25 The latter may be treated as dominant markers (Elshire et al., 2011). Both the co-dominant and
26 dominant tags may be genetically mapped using co-segregation and traditional mapping
27 functions (Elshire et al., 2011; Jiang and Zeng, 1997). This mapping effort is greatly assisted by
28 comparing the DNA sequence of the tags to previously mapped sequences. Tags identified
29 through association studies of large collections of individuals can be used for characterizing
30 germplasm collections or disparate populations (Romay et al., 2013), and for defining marker-
31 trait associations, for example, through a conditional probability approach (Skinner et al., 2000).
32

33 **GIO Features**

34
35 GIO is provided as a flexible tool to carry out the identification of dominant and co-dominant
36 tags and their corresponding alleles within user-specified length, occurrence frequency, and read
37 depth limits; preliminary mapping of codominant markers in bi-parental studies is provided.
38 Association panel studies are supported and result in all of the unique tags found within user-
39
40
41
42
43
44
45
46
47
48

1
2
3 specified length and frequency limits. Chromosomes of origin of the tags are provided in both
4 kinds of studies. The GIO scripts were developed using the PERL scripting language
5 (<http://www.perl.org/>) and numerous tools available as part of the Linux operating system as
6 well as additional open-source tools described below. GIO was developed using NGS results
7 from wheat (*Triticum aestivum* L.), but should be applicable to any organism.
8
9

10 11 12 **Implementation of GIO**

13
14
15 **1) Preprocessing of the data:** GIO provides an optional tag-trimming step in which any bases 5'
16 to a user-specified series of bases are removed and only tags starting with the specified sequence
17 are retained. The intent is to remove bases 5' of the restriction site used to ligate adapters and
18 barcodes in multiplexed sequencing, and to discard any tags that resulted from contamination or
19 sequencing error at the start of the sequencing read. Also provided is an option to delete tags that
20 contain regions of identical bases (homopolymers) of user-specified length. Homopolymeric
21 regions have been shown to be potentially problematic in next-generation sequencing results
22 (Quail et al., 2012).
23
24
25
26
27
28

29
30 **2) Specification of population analysis parameters:** All parameters are specified through a
31 simple graphical user interface window (Fig. 1). Default values are provided and all settings are
32 retained from each run of the scripts, such that user-defined parameters become the default
33 settings for the next run. Basic help information is available from the "Help" button. Data files
34 are assigned a numerical identifier for analysis purposes and for biparental populations, GIO
35 must be told which data files represent the parents. A button is provided to allow the user to view
36 the assigned numerical identifiers, thereby providing the identifiers of the parent files. This
37 button also provides a check that the working directory and file type were specified correctly
38 before starting the analysis. Options available to be specified are: length of tags, minimum
39 frequency among lines, maximum frequency among lines, minimum and maximum number of
40 SNPs per tag, and minimum read depth for a tag to be counted as present. The minimum and
41 maximum frequency among lines parameter is provided to allow the user to screen out very
42 common or very uncommon tags, especially useful when analyzing a segregating population
43 with a 1:1 expected segregation ratio. Once the parameters have been specified, the analysis is
44 initiated by pressing the "Submit" button (Fig. 1).
45
46
47
48
49
50
51
52
53
54
55

56
57 **3) Analysis of a segregating population.**
58
59
60

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
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Co-dominant tags: All unique DNA sequence tags meeting the specified parameters are identified; all tags retained for segregation analysis must occur in only one parent (zero frequency in the second parent), before the minimum read depth parameter is applied. Using SOAP2 (Li et al, 2009), all of these tags are then compared to reference sequence databases. The scripts were developed using a reference database from the International Wheat Genome Sequencing Consortium (IWGSC, 2014) including over 25,000,000 mapped sequences from all of the chromosome arms of wheat (*Triticum aestivum* L.), over 190,000 unique transcript sequences compiled by the first author (http://cereals.bioinfo.wsu.edu/wa_wheat_transcriptome), DNA sequences from the TREP database of transposable elements (<http://wheat.pw.usda.gov/ITMI/Repeats>), as well as chloroplast, mitochondria, and, ribosomal genes sequences. Tags homologous to the organellar, ribosomal, and transposon sequences, within a user-specified number of mismatches, are discarded. The putative chromosome arm of origin is then assigned to all remaining tags if significant homology to the reference sequences was found. If no significant homology was found, the tags are identified as being of “unknown” origin.

Next, of the tags that were retained, those that occurred in the first parent but not the second are compared to all tags from the second parent to identify pairs of tags that differ by the user-specified number of SNP(s), thereby identifying co-dominant markers. The presence/absence of pairs of tags identified in this way are then scored across all progeny lines and genotypes are recorded as homozygous for the first parent genotype, as homozygous for the second parent genotype, as heterozygous, or as unknown if neither parent tag was found.

Once the genotypes of the parents and progeny are identified, the tags are separated into two files. The first file holds “repetitive” tags; these are groups of tags in which one tag from one parent paired with multiple tags from the second parent, differing by the specified number of SNP(s). The second file is identified as holding “nonrepetitive” tags, pairs of tags in which a tag from one parent matched a single tag from the second parent, differing by the specified number of SNP(s). A table of read depth in the parents and all progeny is generated for each of the nonrepetitive tags. The nonrepetitive tags are then genetically mapped in two ways as follows.

De novo construction of genetic maps of co-dominant tags: Usually, several thousand co-dominant, nonrepetitive tags are found. Attempting to genetically map the entire dataset is time-

1
2
3 consuming and requires human intervention to generate meaningful linkage groups. To provide
4 a starting point for these linkage groups, GIO first utilizes clustering algorithms included in R (R
5 Core Team, 2014) to cluster the markers on the basis of cosegregation. The user provides an
6 estimate of the number of clusters that may be found, usually the number of chromosomes or
7 chromosome arms in the organism; R uses this estimate as the clustering basis. By default, GIO
8 uses Euclidean distance measures and clustering using Ward's method. GIO then reformats the
9 segregation data from each cluster to MSTMap (Wu et al., 2008) input format and executes
10 MSTMap to attempt to generate map(s) for each cluster using the Haldane mapping function.
11 Three files are saved from each cluster; the original segregation data, a list of tag pairs and the
12 map position within the cluster, a list of tag pairs that were included in the cluster on the basis of
13 segregation, but were not genetically linked to the other tags in the cluster according to the
14 MSTMap results. A combined file with linkage group and map position of all mapped markers
15 also is provided to facilitate map image creation in other software.
16

17
18 Next, GIO compiles a dataset for each index that was provided to SOAP2 to search for
19 homology to known genomic sequences. In our studies of the wheat genome, each chromosome
20 arm is represented in a separate index so GIO builds separate datasets containing only tags with
21 significant homology to one chromosome arm and tags with no significant homology to
22 previously-mapped sequences. MSTMap is again called to generate map(s) from each of these
23 data sets. From the output, GIO combines linkage groups that contain at least one previously-
24 mapped tag into a single file to facilitate inspection of the linkage groups relative to each other.
25

26
27 ***Dominant tags:*** Tags found in only one of the parents with the user-specified read depth and
28 frequency among the progeny lines, and for which no tags differing by the specified number of
29 SNPs were found in the other parent, are considered to be dominant markers. GIO provides two
30 output files, one with the sequences, chromosome assignments and occurrence among progeny of
31 the dominant markers from the first parent, designated "Parent A", and a file with this
32 information for the dominant markers from the second parent, designated "Parent B." It may be
33 possible to include some of these tags on the co-dominant marker-based maps on the basis of co-
34 segregation. To facilitate this kind of mapping, a file of "neighbor tags" is generated.
35

36
37 ***Neighbor tags:*** Ideally, genetic maps would contain map locations for all available tags, both co-
38 dominant and dominant. However, in our experience with wheat, the great majority of pairs of
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 dominant and co-dominant tags do not co-occur in more than 50% of the progeny lines, and
4 therefore show no evidence of being genetically linked. To facilitate identification and mapping
5 of those markers that may be linked, GIO provides two files of “neighbor” tags, those that co-
6 occur in at least 65% of the lines in the data set. The first, “MarkerNeighbors.txt,” holds the
7 genotypes of the markers in each of the lines designated as A, B, H or U for the co-dominant
8 markers with genotypes of the first parent, second parent, both parents (heterozygous), or neither
9 parent (unknown), respectively. The genotypes of the dominant markers are designated as C, D,
10 or U, indicating the dominant marker from the first parent, second parent, or neither parent,
11 respectively.
12
13

14
15
16
17
18
19
20 **Progeny-only tags:** In our experience with biparental, segregating populations, the majority of
21 the unique tags found occur in the specified frequency of sample lines, but are not found in either
22 parent. Clearly, these tags must occur in one or both parents, but were simply missed by the
23 sequencing process. Optionally, GIO will tally the occurrence of these tags among the sample
24 lines and generate a separate output file with the presence-absence data.
25
26
27
28

29 **4) Analysis of an association panel.**

30
31 Association panels may consist of a number (usually hundreds) of lines that may or may not be
32 related. From a GBS perspective, the objective of association panel analysis is to identify DNA
33 tags that are associated with trait(s) of interest. GIO provides the sequences, chromosome
34 assignments, and frequency of occurrence of all of the tags found in the data set that meet the
35 user-specified criteria of length, frequency and read depth. In addition, a separate data file listing
36 tags that differ by the user-specified number of SNPs is produced. However, in our experience
37 with association panels of hexaploid wheat, a given tag may differ from dozens of other tags by
38 the allowed number of SNPs, all occurring with the allowed frequency among the plant lines.
39 These differences probably are valid, but greatly complicate defining alternate alleles.
40 Therefore, GIO provides a separate output file with pairs of tags that occur with only one “mate”
41 in the data set, i.e. tags that differ by the allowed number of SNPs from exactly one other tag.
42 Identifying the one-mate tags requires the tool agrep (<https://github.com/Wikinaut/agrep>) to be
43 installed. These pairs of tags potentially represent alternate alleles at a single locus.
44
45
46
47
48
49
50
51
52
53

54 **Example Results**

55 **Biparental population**

1
2
3 To test GIO analysis of a segregating population, a dataset comprised of short reads generated by
4 the Ion Proton sequencing system from the parents and 186 progeny from a cross of two winter
5 wheat (*Triticum aestivum*, 2N=6X=42) cultivars was analyzed. There were a total of
6 90,137,056 raw sequence reads. The options specified to GIO were: Length of tags: 75,
7 Minimum frequency among lines: 0.25, Maximum frequency among lines: 0.75, minimum
8 number of SNPs per tag: 1, maximum number of SNPs per tag: 2, minimum read depth: 1, initial
9 estimate of cosegregation clusters: 30. Any tags with homopolymers of five or more consecutive
10 bases were deleted.

11
12 **Co-dominant markers.** The html summary report generated by GIO is shown in Table 1. Of
13 the 19,837 pairs of tags with 1 or 2 SNPs, 4,326 were nonrepetitive, i.e. pairs of tags that differed
14 from each other by one or two SNPs and by more than two from all other tags found. Of these
15 4,326 tag pairs, 289 were homologous to known genes. Working with these 4,326 tag pairs (each
16 tag pair comprises a “marker”) and the clusters found on the basis of cosegregation, GIO used
17 MSTMap to find 58 linkage groups containing 3,354 markers. Next, using the datasets based on
18 homology to previously-mapped markers, GIO found 2,564 markers in 428 linkage groups
19 (about 10 linkage groups per chromosome arm), where each linkage group contained at least one
20 marker homologous to a previously-mapped marker. These 2,564 markers were mapped to
21 8,282 locations, meaning that each marker mapped to an average of 3.2 chromosome arms.
22 Because this mapping was in hexaploid wheat, a map location on three chromosome arms, one in
23 each of the three diploid genomes, is not surprising. An example of the linkage groups found by
24 GIO from wheat chromosome 1A is shown in Supplementary Figure S1.

25
26 **Dominant tags.** There were 45,982 apparently dominant tags found from the first parent, and
27 27,479 from the second. In the comparison of dominant tags to the 2,363 nonrepetitive
28 codominant markers (reported by GIO as “neighbor” tags), 1,441,909 combinations were found
29 that occurred simultaneously in at least 65% of the plant lines. These combinations involved all
30 2,363 co-dominant markers and 4,419 dominant tags. Most co-dominant markers were associated
31 with a relatively small number of dominant tags. The number ranged from one to 2,993; the first
32 quartile was 18, the median was 220, and the third quartile was 1,041.

33
34 Similarly, the dominant “neighbor” tags were each associated with relatively few co-dominant
35 markers; the first quartile was 13, the median was 195, the third quartile was 587, and the
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 maximum was 1,319. While the segregation data of all of the dominant tags are provided in the
4 GIO output (files from the two parents are named ParentAdominantSNPS.txt and
5 ParentBdominantSNPS.txt, respectively), our experience has shown that generating maps from
6 the raw segregation data is essentially impossible with the mapping software we have tried.
7 Thus, the “neighbor” tags are provided to help augment the genetic maps based on co-dominant
8 markers; especially useful in regions of interest found through marker-trait association studies.
9 For example, in the biparental mapping population described above, 335 co-dominant markers
10 were mapped to chromosome 1A. Of those 335 markers, 211 were associated with at least one
11 dominant tag in at least 65% of the plant lines.
12
13
14
15
16
17
18

19 **Association Panel**

20 To test GIO analysis of an association panel, a genotyping by sequencing approach was applied
21 to a panel of 402 unique winter wheat entries with 607,325,655 individual sequence reads
22 generated by the Ion Proton system. With the frequencies of accepted tags limited to 5-95% (i.e.
23 excluding the very rare and the very common tags), 611,456 unique tags were found. Of these,
24 there were 54,551 pairs of tags that differed by a single SNP from only one other tag in the
25 dataset. The genotypes of each of the plant lines for these 54,551 tag pairs were reported as
26 homozygous for the first or second tag, as heterozygous, or as unknown. The read depth, i.e. the
27 number of times each tag occurred in each plant line, was determined and reported in a separate
28 file. These results are reported as text files which may be used with other software for
29 determining marker-trait associations.
30
31
32
33
34
35
36
37
38
39

40 **Conclusions**

41 GIO provides a highly customizable, simple-to-use tool that starts with raw sequence read files,
42 cleans up the read data as specified by the user, and generates genotype, frequency and linkage
43 group information with very little user intervention. Segregating populations and association
44 panels are supported, resulting in comprehensive tabulation of tag occurrence data in simple text
45 format, enabling downstream processing by other software to discover marker-trait associations.
46
47
48
49
50
51

52 **Acknowledgements**

53
54
55
56
57
58
59
60

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
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

This project was supported by USDA–ARS project 5348-21000-030-00D. Mention of product names does not represent an endorsement of any product or company but is given only to clarify the methodology; other products may be equally effective.

For Review Only

References

1. Elshire, R.J., Glaubitz, J.C., Sun, Q., Poland, J.A., Kawamoto, K., Buckler, E.S., et al. 2011. A Robust, Simple Genotyping-by-Sequencing (GBS) Approach for High Diversity Species. PLoS ONE. 2011 May 4;6(5):e19379.
2. IWGSC: The International Wheat Genome Sequencing Consortium. 2014. A chromosome-based draft sequence of the hexaploid bread wheat (*Triticum aestivum*) genome. Science, 345(6194). <http://doi.org/10.1126/science.1251788>
3. Jiang, C., and Zeng, Z-B. 1997. Mapping quantitative trait loci with dominant and missing markers in various crosses from two inbred lines. Genetica. 101(1):47–58.
4. Lander, E.S., and Botstein, D.1989. Mapping mendelian factors underlying quantitative traits using RFLP linkage maps. Genetics. 1989 121(1):185–99.
5. Li, R., Yu, C., Li, Y., Lam, T-W., Yiu, S-M, Kristiansen, K., et al. 2009. SOAP2: an improved ultrafast tool for short read alignment. Bioinformatics. Aug 1;25(15):1966–1967.
6. Quail, M.A., Smith, M., Coupland, P., Otto, T.D., Harris, S.R., Connor, T.R., et al.2012. A tale of three next generation sequencing platforms: comparison of Ion Torrent, Pacific Biosciences and Illumina MiSeq sequencers. BMC Genomics. 24;13(1):341.
7. R Core Team. 2014. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/>.
8. Romay, M.C., Millard, M.J., Glaubitz, J.C., Peiffer, J.A., Swarts, K.L., Casstevens, T.M., et al. 2013. Comprehensive genotyping of the USA national maize inbred seed bank. Genome Biology. 11:14(6):R55.
9. Sachidanandam, R., Weissman, D., Schmidt, S.C., Kakol, J.M., Stein, L.D., Marth, G., et al. 2001. A map of human genome sequence variation containing 1.42 million single nucleotide polymorphisms. Nature. 409(6822):928–33.
10. Skinner, D.Z., Loughin, T., Obert, D.E. 2000. Segregation and conditional probability association of molecular markers with traits in autotetraploid alfalfa. Molecular Breeding. 6(3):295–306.

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
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

11. Wu, Y., Bhat, P.R., Close, T.J., Lonardi, S. 2008. Efficient and accurate construction of genetic linkage maps from the minimum spanning tree of a graph. PLoS Genet. 2008 Oct 10;4(10):e1000212.

For Review Only

Figure Caption

Figure 1. The graphical user interface (GUI) used to interact with GIO. All options are specified through this one GUI.

Supplemental Figure S1. Example of the linkage groups found by GIO from wheat chromosome 1A. The genetic distances were computed by MSTMap. The underlined numbers are the genetic positions described in the IWGSC mapping effort of Chinese Spring for markers that were homologous to the markers found by GIO and identified on the map. The chromosome arm designation following a marker name indicates significant homology of the corresponding marker to sequence(s) from the indicated chromosome was found in comparisons to the Chinese Spring data from the IWGCS mapping effort. The images were drawn using MapChart (Voorrips, R.E., 2002. MapChart: Software for the graphical presentation of linkage maps and QTLs. The Journal of Heredity 93 (1): 77-78.)

For Review Only

Table 1. Summary of tag, polymorphism, genotype information, and runtime specifics as output by GIO.

Type of tag	Tag Count	SNP	SNP Frequency	Genotypes among progeny	Genotype Frequency
Initial Unique Tags of 75 bases	18,263,477	[K] (G/T)	8,063	A (1st parent)	927,972
Tags meeting frequency criteria	189,740	[M] (A/C)	9,561	B (2nd parent)	884,747
Tags from 1st parent with SNP(s) relative to 2nd parent	9,405	[R] (A/G)	16,957	H (heterozygous)	534,889
Tags from 2nd parent with SNP(s) relative to 1st parent	7,745	[S] (C/G)	10,645	U (unknown)	1,381,748
Tag pairs with SNP(s)*	19,837	[W] (A/T)	6,477	Codominant	
Dominant tags, 1st parent	45,982	[Y] (C/T)	15,913	Polymorphism	Count
Dominant tags, 2nd parent	27,479			SNPs	18,478
Tags in progeny only**	99,129			INDELS	1,359

*Tags from one parent may pair with more than one tag from the second parent.

**Tags that met the length and frequency criteria but were not found in either parent.

Runtime specifics

The total number of sequence read files found was 188.

The parents were represented in files numbered 1 and 2.

The runtime options were:

Length of tags: 75

Minimum frequency among lines: .25

Maximum frequency among lines: .75

Minimum number of snps per tag: 1

1
2
3
4 Maximum number of snps per tag: 2

5 Minimum number of times a tag had to occur in a sequence file to be included
6 (accepted as "real"):1
7

8
9 The number of independent processes (approximately the number of cores used)
10 was: 30

11 SOAP was used to search for homology to the sequence databases

12 in: **/home/dzs/GBS/Kim/SNPscripts/GIO/Filters** allowing 1 mismatch(es) per
13 tag.
14

15 Tags homologous to sequences in Repeats.fa were discarded.

16 Tags homologous to sequences in wwt.fas were identified as putative genes.

17 Adjacent base transposition was counted as a single SNP.
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Interface to G.I.O. Genes In Order, a collection of PERL and BASH shell scripts for variant DNA sequence tag identification

Welcome to G.I.O.
Genes In Order

CHECK THIS BOX TO TURN ON TAG CLEAN UP AND SPECIFY OPTIONS BEFORE ANALYSIS CHECK THIS BOX FOR ASSOCIATION STUDY (as opposed to segregation)

Directory of sequence reads? File extension of sequence reads?

Press this button to show the sequence read data files and their corresponding number in the specified directory and file type:

First parent is file number: Second parent is file number:

DNA Tag specs:

Length of tags? Minimum times a tag sequence must occur in a sample to be included?

Minimum frequency among progeny or accession lines? Maximum frequency among progeny or accession lines?

Minimum number of variant bases per tag? Maximum number of variant bases per tag?

How many computer cores are available for this job? About how many linkage groups do you think are present?

Do homology searches with the SOAP indices in this directory: Allowing this many mismatches per tag:

Omit unwanted tags (e.g. organelles) homologous to sequences in this database:

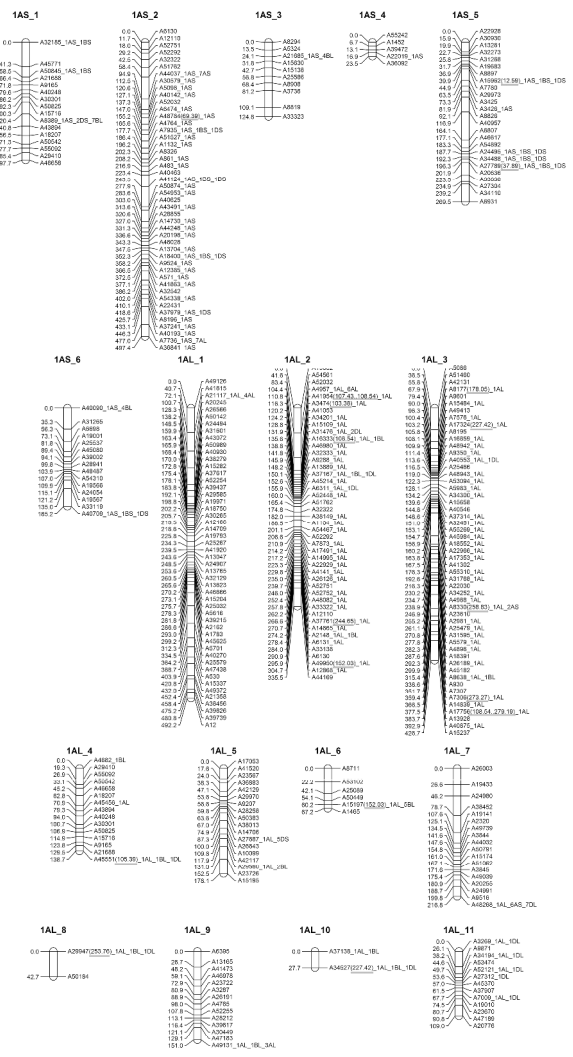
Label tags homologous to sequences in this database as genes:

Count adjacent base transposition as one change Tally occurrence of tags found only in the progeny

Figure 1. The graphical user interface (GUI) used to interact with GIO. All options are specified through this one GUI.

479x270mm (300 x 300 DPI)

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
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



Supplemental Figure S1. Example of the linkage groups found by G10 from wheat chromosome 1A. The genetic distances were computed by MSTMAP as part of the G10 process. The underlined numbers are the genetic positions described in the IWGSC mapping effort of Chinese Spring for markers that were homologous to the markers found by G10 and identified on the map. The chromosome arm designation (following a marker name indicates significant homology of the corresponding marker to sequences) from the indicated chromosome was found in comparisons to the Chinese Spring data from the IWGSC mapping effort. The images were drawn using MapChart (Voorrips, R.E., 2002. MapChart: Software for the graphical presentation of linkage maps and QTLs. The Journal of Heredity 93 (1): 77-78.)

304x632mm (300 x 300 DPI)