

This provisional PDF was built from the peer-reviewed and accepted manuscript submitted by the author(s).
The manuscript has not been copyedited, formatted or proofread.
Please note that the provisional version can differ from the final version.
Final fully formatted version will be available soon.

Short communication

A database for efficient storage and management of multi panel SNP data

Eildert Groeneveld and Cong VC Truong

Department of Breeding and Genetic Resources, Institute of Farm Animal Genetics (FLI), Neustadt, Germany

For information about "Archiv Tierzucht" please visit <http://www.archivtierzucht.de/>.

Archiv Tierzucht 56 (2013) 103
doi: 10.7482/0003-9438-56-103

Received: 17 July 2013
Accepted: 19 November 2013
Online: 20 November 2013

Corresponding author:

Eildert Groeneveld; email: eildert.groeneveld@fli.bund.de
Department of Breeding and Genetic Resources, Institute of Farm Animal Genetics (FLI), Neustadt, Germany

© **2013 by the authors**; licensee Leibniz Institute for Farm Animal Biology (FBN), Dummerstorf, Germany.
This is an Open Access article distributed under the terms and conditions of the Creative Commons Attribution 3.0 License (<http://creativecommons.org/licenses/by/3.0/>).

1 **A database for efficient storage and**
2 **management of multi panel SNP data**

3 Eildert Groeneveld*, Cong VC Truong

4 **Department of Breeding and Genetic Resources**

5 **Institute of Farm Animal Genetics (FLI)**

6 **D-31535 Neustadt, Germany**

7 ***E-mail: eildert.groeneveld@fli.bund.de**

8 **Abstract**

9 The fast development of high throughput genotyping has opened up new possi-
10 bilities in genetics while at the same time producing immense data handling is-
11 sues. A system design and proof of concept implementation are presented which
12 provides efficient data storage and manipulation of single nucleotide polymor-
13 phism (SNP) genotypes in a relational database. A new strategy using SNP
14 and individual selection vectors allows us to view SNP data as matrices or sets.
15 These genotype sets provide an easy way to handle original and derived data,
16 the latter at basically no storage costs. Due to its vector based database storage,
17 data imports and exports are much faster than those of other SNP databases.
18 In the proof of concept implementation, the compressed storage scheme reduces
19 disk space requirements by a factor of around 300. Further, this design scales
20 linearly with number of individuals and SNPs involved. The procedure supports
21 panels of different sizes. This allows a straight forward management of different
22 panel sizes in the same population as occur in animal breeding programs when
23 higher density panels replace previous lower density versions.

24 **Introduction**

25 High throughput single nucleotide polymorphism (SNP) genotyping is evolving
26 at a staggering rate, developing into a powerful tool in genetic analyses in all
27 areas of biology. While its promises are immense, so are the data processing
28 issues associated with it. Dropping genotyping costs and ever increasing marker
29 densities result in a huge increase in data volume, which seems to develop faster
30 than the already impressive rate at which data storage costs have come down
31 in the past. Thus, increasing data storage requirements are an issue.

32 SNP data analysis workflows often result in filtering SNPs thereby creating
33 genotype subsets for different purposes. This process can lead to a dramatic
34 increase in disk space requirement: while each derived dataset will always be
35 smaller than the original, their sizes will still be of the same magnitude, quickly
36 using huge amounts of disk space.

37 High throughput genotyping is used across species with diverse sets of geno-
38 typing panels of different sizes ranging from a few thousand to millions. New
39 panels of higher densities may be used to retype the same individuals as is done
40 in regular animal breeding programs, resulting in ever growing datasets. Ac-

41 cordingly, information originating from panels of different densities will have to
42 be managed and processed.

43 Space requirements of SNP data may be very large, resulting in specialized
44 hardware having to be made available for storage. Therefore, it makes sense to
45 centralize data management in a relational database, which can be accessed by
46 multiple users. A number of database developments try to address some of the
47 above issues [3, 1, 2].

48 However, the main shortcoming is the 'one genotype per row' (OGPR) stor-
49 age scheme which leads to huge storage requirements and slow export times.
50 Rios et al.(2010) improved on this scheme by storing one record per individ-
51 ual containing all genotypes together with each genotype's position. Our pro-
52 posal goes well beyond this through a more efficient storage scheme and the
53 development of the genotype set concept, which substantially enhances data
54 manipulation.

55 Efficient management of SNP data has to address long term data storage,
56 multiple genotyping panels of different sizes as are already in practical use in
57 animal breeding and elsewhere. Further, efficient selection of SNPs and fast
58 exports to various formats for further processing are essential. This leads us to
59 the following design.

60 **The Design**

61 **Compressed storage**

62 Typical files from genotyping labs easily have sizes of hundreds of Mbyte per
63 individual. We propose to store the SNP genotypes of one individual in a
64 compressed vector using the position as determined by the panel map, which
65 leads to one genotype record per individual. When only the biallelic state of a
66 SNP is of interest, two bit storage is sufficient, allowing 16 SNPs to be stored
67 in one 32bit integer word. Accordingly, all genotypes from a 60,000 SNP panel
68 can be stored in an integer vector of dimension 3,750 or 14.6 kb. Once a panel
69 map is stored, any number of resulting genotypes can be loaded.

70 Often the call rates i.e. the GCscore need to be stored for each SNP, typically
71 a floating point number in the range of 1 to 100. Here, a flexible scheme is
72 proposed, allowing the user to determine the number of bits to be used for the
73 score: 4 bits accommodate a resolution of 100/15. Using the same 4 bits on the
74 range of 51% to 100% will double the resolution. As with the genotypes, the
75 GCscore is also designed as one vector per individual.

76 **Set based manipulation**

77 For easy data manipulation, we introduce the concept of effectively spaceless
78 genotype sets. Each SNP panel comprises a specified set of SNPs. Using the
79 SNP's position in the panel as its position in the genotype bit vector enables
80 access to each SNP without explicitly having to state the SNP name. Once

81 genotypes are treated as vectors, a SNP dataset can be viewed as a matrix
82 of genotypes with the SNPs constituting the columns while each genotyped
83 individual leads to one row.

84 Often, only subsets of SNPs are used in analyses, perhaps only those from a
85 particular chromosome, or SNPs with a minimum frequency. Using a bit vector
86 *snp_sel_vec* of the dimension of the SNP panel provides a generalized approach:
87 a .TRUE. or .FALSE. decides if a SNP is a member of a particular subset. A
88 genotype set is then completely defined by adding a vector *indiv_sel_vec* which
89 contains the individuals of the subset. Finally, a set name for the combination
90 of a particular *snp_sel_vec* and *indiv_sel_vec*, uniquely specifies a genotype set.
91 Thus, creating new derived genotype sets amounts to creating two new lists: one
92 for the SNPs and the other for the individuals and giving this a genotype set
93 name for easy access. The storage implications are clear: instead of having to
94 store a matrix of dimension $nSNP * nIndividual$ only two vectors of dimension
95 $nSNP$ and $nIndividual$ need to be stored for each derived matrix which are
96 thus effectively spaceless.

97 It should be noted that the concept of genotype sets allows rapid implemen-
98 tation of general set operations like unions or intersections on the basis of any
99 genotype set, original or derived.

100 Being a central repository for all data to be analysed, fast exports are critical.
101 A data analysis step typically starts with an export using an appropriate format
102 for downstream processing. Often exports are performed through costly SQL
103 selects on the basis of SNP names [1, 2, 3].

104 However, export speed is a function of both the data storage and the retrieval
105 scheme. Storage overhead of our approach is minimal: the *indiv_sel_vec* is an
106 integer vector with as many 32bit words as there are individuals in the genotype
107 sample, while *snp_sel_vec* has the dimension of the number of SNPs in the panel
108 divided by 16 (for 2 bit storage). Thus, all genotypes from a 1 million panel
109 will occupy 244kb in *snp_sel_vec*. Clearly, on this data volume basis our design
110 will have a performance advantage over the OGPR paradigm which needs to
111 process 1 million records per individual.

112 For data retrieval, the genotype set approach replaces SQL based SNP selec-
113 tion by much faster vector operation. An export of a genotype set amounts to
114 these actions: firstly, *snp_sel_vec* and *indiv_sel_vec* are fetched for the chosen
115 set. Secondly, for each individual the compressed genotype vector is retrieved
116 through one SQL select and shrunk on the basis of the *snp_sel_vec* which can
117 be implemented as fast shifts. Thus, the extraction speed is largely independent
118 from the subset selection.

119 Performance of Proposed Design

120 The proof of concept implementation was done in Perl and PostgreSQL as the
121 backend database server. Apart from conceptual simplicity, performance in
122 terms of storage efficiency and speed of data import and export is of critical
123 importance. This was investigated through a number of genotype datasets of

124 very different sizes (Table 1).

125 The timings refer to the import of SNPs after the panel map has been loaded,
126 as this is done only once. Further, two bit storage for the biallelic state genotypes
127 (A, B, AB, and no call) with no GCscore is assumed. The benchmarks were
128 executed on an iCore 5 laptop (2.7GHz / 4GB RAM).

129 Importing data following our design, requires two steps. Firstly, the panel
130 map is loaded which contains the panel size, the SNP names and chromosome
131 location. A unique panel name is given, which needs to be specified when-
132 ever new genotype data is loaded in a second step. During the initial load, a
133 *snp_sel_vec* is created with all bits set to 1. The individual IDs are stored in
134 *indiv_sel_vec* as picked up from the data file. The combination of those two
135 selection vectors then yields the first genotype set, uniquely identified by its
136 symbolic name, which is the only information required for an export. Notice,
137 that our design scales very well (column 5 and 6 in Table 1, with import and
138 export time per 1 million SNP even going down as panel sizes increase.

139 Popular SNP database management systems store OGPR [1, 2, 3] resulting
140 in huge storage requirements: using the HumapHap300 for 300 subjects, about
141 90 million of records will be generated [2]. In contrast, our design will create only
142 300 rows, which will clearly reduce both storage and processing requirements.
143 Even the improved design presented by Rios et al.(2010) requires 5GB for 1.5
144 billion genotypes, while our design would require approximately .37GB.

145 A direct comparison to other implementations for computing time is difficult
146 to make. Based on the timings from our design (1.55sec and .54/1 million geno-
147 types, for import and exports as a weighted mean from Table 1), we expanded
148 the software comparisons given in Table 1 in Mitha et al. (2011) for a 317K
149 panel with 100 samples. The timings in minutes are n/a, 18, 4.2 and 0.82, for
150 imports and 10, 93, 5, 0.29 for exports from SNPLims, GWASA, SNPpy single,
151 and our implementation, respectively. Thus, our design seems to be an order of
152 magnitude faster than the best of the contenders. This is not surprising, since
153 the OGPR requires space for the sample ID, the SNP name and the chromo-
154 some location for each genotype. Using the database design given in Mitha et al.
155 (2011) for the 317K dataset, we have a storage requirement of 2,721MB versus
156 8MB. Thus, our compressed bit vector storage scheme is around 300 times more
157 efficient than the OGPR scheme of other packages.

158 As an example, on a 500GB disk 2 million individuals genotyped with a
159 317K panel can be stored using our scheme. In contrast, the same disk would
160 store .146 million individuals with the design from Rios et al.(2010) , and only
161 7000 with the common OGPR storage scheme.

162 Supplementary Material

163 The proof of concept code is available under the GPL license at
164 <ftp://ftp.tzv.fal.de/pub/snp/SNP-PoC.tar.gz>

Table 1. Performance for five test datasets

data	n SNP ^a	n ind ^b	tot SNP ^c	imp ^d	exp ^e	DB stor ^f
set1	58	47	2.7	2.38	1.38	0.28
set2	36	403	14.7	4.37	1.02	0.27
set3	229	90	20.6	1.89	0.69	0.25
set4	4098	270	1106	1.60	0.53	0.25
set5	1458	1397	2036	1.50	0.54	0.25

^apanel size: number of SNPs * 1,000

^bnumber of individuals

^ctotal number of SNPs in dataset * 1,000,000

^dtime to import 1 mio SNPs in seconds

^etime to export 1 mio SNPs in seconds

^fstorage in MByte per 1 mio SNPs

165 References

- 166 1. C. Fong, D. C. Ko, M. Wasnick, M. Radey, S. I. Miller, and M. Brit-
167 tnacher. GWAS analyzer: integrating genotype, phenotype and public
168 annotation data for genome-wide association study analysis. *Bioinformat-*
169 *ics*, 26(4):560–564, January 2010.
- 170 2. F. Mitha, H. Herodotou, N. Borisov, C. Jiang, J. Yoder, and K. Owzar.
171 SNPpy-Database management for SNP data from GWAS studies. *Duke*
172 *Biostatistics and Bioinformatics (B&B) Working Paper Series*, page 14,
173 2011.
- 174 3. A. Orro, G. Guffanti, E. Salvi, F. Macchiardi, and L. Milanese. SNPLims:
175 a data management system for genome wide association studies. *BMC*
176 *Bioinformatics*, 9(Suppl 2):S13, 2008.
- 177 4. Daniel Rios, William M. McLaren, Yuan Chen, Ewan Birney, Arne
178 Stabenau, Paul Flicek, and Fiona Cunningham. A database and API for
179 variation, dense genotyping and resequencing data. *BMC bioinformatics*,
180 11(1):238, 2010.