# Edinburgh Research Explorer

## Locating a novel autosomal recessive genetic condition using only WGS data from three cases and six controls; a case study of a new variant in the cattle glucokinase gene

OPEN ACCESS

# Locating a novel autosomal recessive genetic disease using only WGS data from three cases and six controls;  a case study of a new variant in the cattle glucokinase gene

1    **Geoffrey Pollott[1*], Richard Piercy[2], Claire Massey[2], Mazdak Salavati[3,4] and D. Claire Wathes[3]**

2    [1]Department of Pathobiology and Population Sciences, The Royal Veterinary College, Royal College
3    Street, London, NW1 0TU, UK

4    [2]Comparative Neuromuscular Diseases Laboratory, Department of Clinical Sciences and Services,
5    Royal Veterinary College, Royal College Street, London NW1 0TU, UK

6    [3]Department of Pathobiology and Population Sciences, The Royal Veterinary College, Hawkshead
7    Lane, Hatfield, Hertfordshire, AL9 7TA, UK

8    [4]Now at The Roslin Institute, Easter Bush, Midlothian, EH25 9RG, Scotland, UK

9    **\* Correspondence:**
10   Corresponding Author
11   gpollott@rvc.ac.uk

13   **Abstract**

14   New Mendelian genetic conditions arise all the time. To manage them effectively methods need to be
15   devised that are quick and accurate. Until recently, finding the causal genomic site of a new
16   autosomal recessive genetic disease has required a two-stage approach using single nucleotide
17   polymorphism (SNP) chip genotyping to locate the region containing the new variant. This region is
18   then explored using fine-mapping methods to locate the actual site of the new variant. This paper
19   explores bioinformatic methods that can be used with just nine whole-genome sequenced animals to
20   simplify and expedite the process to a one-step procedure.

21   Using whole genome sequencing of only three cases and six controls, the site of a novel variant
22   causing perinatal mortality in Irish Moiled calves was located. Four methods were used to interrogate
23   the variant call format (VCF) datafile of these nine animals; genotype criteria (GCR), autozygosity
24   by difference (ABD), variant prediction scoring and registered SNP information. From more than 8
25   million variants in the VCF file only one site was identified by all four methods (Chr4:
26   g.77173487A>T (ARS-UCD1.2 (GCF_002263795.1)). This site was a splice acceptor variant located
27   in the glucokinase gene (*GCK*). It was verified on an independent sample of animals from the breed
28   using genotyping by polymerase chain reaction at the candidate site and autozygosity by difference
29   using SNP chips. Both methods confirmed the candidate site. Investigation of the GCR method found
30   that sites meeting the GCR were not evenly spread across the genome but concentrated in regions of
31   long runs of homozygosity. Locating GCR sites was best done using two controls to every case and
32   the controls should be distantly related to the cases, within the breed concerned. Less than 20 animals
33   need to be sequenced when using the GCR and ABD methods together.

34 The genomic site of novel autosomal recessive Mendelian genetic diseases can be located using less
35 than 20 animals combined with two bioinformatic methods, autozygosity by difference and genotype
36 criteria. In many instances it may also be confirmed with variant prediction scoring. This should
37 speed-up and simplify the management of new genetic diseases to a single-step process.

38 **1    Introduction**

39 A review of 34 papers detailing work to map 38 novel autosomal recessive genetic conditions to their
40 position on the genome (Pollott, 2018) suggested that finding the site of such a condition required the
41 use of at least a two-stage methodology. Firstly the use of a suitable single nucleotide polymorphism
42 (SNP) chip in a case/control study to locate the *region* containing the new variant combined with a
43 method searching for long runs of homozygosity (ROH), the more traditional chi-squared method
44 being shown to be inadequate. The second stage used a range of 'fine-mapping' methods to search
45 within the highlighted region for the *site* of the new variant, many of which resulted in whole genome
46 sequencing (WGS) of a few cases and controls. More recently methods have been developed which
47 use WGS as an initial step but such methods typically require additional resources or sequencing a
48 large number of animals.

49 The objective of the current paper was to see if it is possible to locate the position of a novel
50 autosomal recessive genetic condition directly using the ideas contained in Pollott (2018) on WGS
51 methodology and using a small number of cases and controls without recourse to more extensive
52 resources which may not always be available. To achieve this, we investigated combining this
53 approach with a range of other bioinformatic tools and genetic ideas which may indicate the site of
54 such a variant by reading the various signals in the WGS data. Using the Variant Call Format (VCF;
55 VCF (2019)) files from a suitable combination of cases and controls it was suggested that typically
56 about 16 animals would need to be sequenced in order to locate a new autosomal recessive variant
57 using a 'genotype criteria' approach (GCR; Pollott (2018)).

58 Here we test these ideas on a novel genetic disease found in Irish Moiled cattle. The Irish Moiled is
59 an ancient hornless breed native to the island of Ireland. It was popular in the 1800s but by the late
60 1970s the pedigree herd numbered only 30 breeding females and two bulls. In 1979 the Rare Breeds
61 Survival Trust recognised the Irish Moiled as endangered and placed the breed on its 'critical' list
62 (Irish Moiled Cattle Society, 2020). The population size now numbers about 875 females and 90
63 bulls. Fortunately, novel fatal genetic diseases are relatively rare but when they do occur it is
64 important to find the cause and implement plans to manage the condition via selective breeding, as
65 soon as possible. A number of Irish Moiled cattle breeders were concerned about the seemingly high
66 occurrence of early calf deaths in their herds. Affected calves (of both sexes) appeared normal at
67 birth, and were initially active, but then deteriorated rapidly, dying within the first week after birth.
68 An initial analysis was undertaken which suggested that there was likely a genetic basis to the disease
69 (see Supplementary File Page S18) and, since it was fatal, that it could only be inherited as a
70 recessive condition.

71 **2    Materials and Methods**

72 Throughout this paper the ARS-UCD1.2 (GCF_002263795.1) build of the cattle genome was used
73 and all chromosome positions quoted relate to it (Ensembl, 2020).

74 **2.1    Sample Collection**

**Commented [MS1]:** We should really clarify how the coordinates of variants from UMD3.1 were lifted over to ARS-UCD1.2.
I suggest you add this as a step 9 in the supplementary methods section S2-S3.

Without this clarification it looks confusing for the reviewer as why the data wasn't originally mapped to ARS-UCD1.2.

**Commented [MS2]:** Its beneficial to mention the version of Ensembl. Current version is v104 but I'd assume you have used either 102 or 103 given the timeline of your work.

75　Farmers were contacted via the Irish Moiled breed society and asked to submit hair samples for
76　analysis. Clean tufts of tail hairs were plucked from either live cows/bulls or recently dead calves
77　(within 8 h post mortem) by the animals' owners. These were sealed in a paper envelope for posting
78　to the laboratory. Information recorded included: sample type (bull, dam or dead calf); ear tag
79　numbers of dam and sire; dam herd ID; calving date; calf sex; time of calf death (stillborn/ days after
80　calving). DNA was then extracted from the hair follicles for processing as described in the
81　Supplementary File (Page S2).

82　**2.2　Whole genome sequencing**

83　Nine DNA samples were used in these analyses comprising three dead calves (cases) and six controls
84　which were either parents of the cases or parents of other dead calves (and grandparents of the cases)
85　as illustrated in the pedigree (Supplementary File Figure S3). These nine samples were sequenced on
86　an Illumina NGS platform after sample preparation as described in the Supplementary File (Page S2).
87　Briefly, 1 μg of DNA per sample was processed using a TruSeq Nano DNA LT Library Prep Kit
88　(Illumina, USA) according to the supplied protocol. This produced randomly sheared 350bp inserts.
89　After end repair and adaptor ligation, DNA was amplified via polymerase chain reaction (PCR) and
90　the product purified using AMPure XP (Beckman Coulter, UK). Size-selected DNA from each
91　animal was sequenced on the HiSeq machine to achieve 150bp paired-end reads to cover the bovine
92　genome with an average 30X coverage (> 90 Gbp raw data with > 85% Q30 (Phred scaled)).
93　Alignment, mapping, variant calling and preparation of the final VCF file were carried out on the
94　subsequent reads as described in the Supplementary File (Page S2).

95　**2.3　Genotype criteria**

96　Information derived from WGS data on a small sample of cases and parental controls may contain a
97　number of signals indicating the site of a novel autosomal recessive condition. 'Across'-animal data
98　should show a typical pattern of homozygous cases and heterozygous parental controls at the
99　candidate site; the 'genotype criteria' approach (Pollott, 2018). Considering a single base position
100　with a reference allele A for the given species and a new variant C which causes a novel autosomal
101　recessive genetic disease, then the expected outcomes from matings between carriers in the
102　population will be offspring with the genotypes AA, AC and CC in the classical Mendelian ratio of
103　1:2:1. Lethals (CC) would be observed in the dataset if the effect of the new homozygous variant
104　occurred after recording. The term 'genotype criteria' (GCR) was used to mean the particular
105　combination of case and control genotypes required to indicate that a base position could harbour the
106　novel lethal variant (Pollott, 2018). For example, in a dataset comprising five cases and 10 parental
107　controls we would expect to find the novel lethal variant at a position showing CC genotypes in all
108　five cases and a genotype containing the C allele in the 10 parental controls, or all AC in the case of a
109　biallelic position. The probability of occurrence of the GCR for this condition would be $1/3^n$ in cases
110　and $1/3^m$ in parental controls, where n is the number of cases and m is the number of controls that
111　have been whole genome sequenced (Pollott,2018). If the VCF datafile comprised 14 million
112　positions (~0.005 of the cattle genome) then a minimum of 15 animals would probably need to be
113　genotyped in order to find one position with the required genotype criteria i.e. $1/3^{15}$ x 14 million =
114　0.98 (i.e. ~1), the expected number of sites with the 'correct' genotype criteria from the genome of 15
115　animals.

116　A script was written in Perl 5.28 to scan the final VCF file (containing all nine WGS animals) for the
117　expected GCR pattern across cases and controls (i.e. all cases homozygous for the same allele and all

**Commented [MS3]:** Is it possible to include the code as a supplementary file?

3

118    controls heterozygous and containing this allele). In order to qualify for selection a site had to have
119    all genotypes with a Phred-scaled quality score greater than 12 and a depth of coverage more than 11
120    reads. The identity of these sites was stored along with their relevant VCF record for later scanning
121    and use.

122    **2.4      Autozygosity by difference and runs of homozygosity**

123    'Within'-animal data should show long ROH around the new variant in cases, which are not present
124    in controls. Variants causing a novel autosomal recessive genetic disease are expected to carry with
125    them a very long haplotype originating from the animal in which the variant first arose. When a new
126    case animal is formed then it contains two copies of this long haplotype, only broken up by any
127    recombination events that have occurred since the formation of the original variant, and the new
128    variant will be situated in a long ROH. This idea has been the basis for locating novel variants using
129    SNP chips for a number of years (see Pollott (2018) for a review) and can be used in a number of
130    ways with VCF data. Long ROH throughout the genome could be found and one would expect to
131    find the novel variant in the longest ROH in cases, possibly with adjustment for the situation in
132    controls. Figure 1 shows a hypothetical section of a chromosome containing a novel autosomal
133    recessive variant for a number of cases and controls to illustrate both the ABD and GCR methods.

134    The autozygosity-by-difference (ABD) method measures runs of homozygosity on each chromosome
135    of each individual in the dataset, cases and controls, using genomewide single-nucleotide variant
136    (SNV) (or SNP) genotypes (Pollott, 2018). Mean ROH length, in Kb, at each SNV positon is
137    calculated for cases and controls separately and then their difference calculated as the ABD score.
138    The likely site of the new variant is in the region with the greatest mean ROH in cases, after taking
139    into account any breed-specific ROH found in controls i.e. the ABD score. The ABD method was
140    programmed in Perl 5.28. The final VCF file was used as the basis to generate a file of sites as input
141    to the ABD method. This method is sensitive to incorrectly called genotypes and so the VCF file was
142    subject to hard filtering as recommended by the Broad Institute (2020) in the absence of suitable
143    databases to use for the recommended variant quality score recalibration (VQSR). A file of SNV
144    were generated from the final VCF file which passed the quality control tests shown in Table S1
145    following a summary of the quality statistics of the VCF file (see Figures S1 and S2).

146    The ABD scores were used to look for the potential site of the novel variant causing calf-mortality in
147    the Irish Moiled dataset. The probability of each ABD score was tested using 100 permutations of the
148    data based on the random allocation of animals to phenotypes and recalculation of the ABD scores
149    (Pollott, 2018). Significance at the $P < 0.01$ level was considered as an indicator of a possible site of
150    the new variant.

151    **2.5      Sorting intolerant from tolerant (SIFT) score**

152    In a fatal genetic disease one would expect to find that the products from any change in the sequence
153    would have a drastic effect on the phenotype of the animal so one could not only search for SNV
154    with a potentially drastic effect but could also eliminate those with 'silent' changes. The Variant
155    Effect Predictor (VEP; McLaren et al., 2016)) is a bioinformatic tool which can take a change in a
156    base at a given position on the genome and predict the outcome of that change on the corresponding
157    coding or non-coding genomic feature. Using this method it is possible to model each base position
158    in the VCF file to see the effect of the new variant on the phenotype in the form of a SIFT score.

159    The final VCF file from all the animals was annotated for variant effect prediction using ~~ENSEMBl~~
160    Ensembl VEPtools command line v90.5 (McLaren et al., 2016) given the following flags: --tab --fork

4

161    8 --offline --species bos_taurus_merged. The VEP was used on all variant sites in the merged VCF
162    file and the results filtered for HIGH SIFT scores using VCFtools (Danecek, 2011). Various
163    outcomes are given in the VEP but here the HIGH outcome was used for the SIFT (Ng and Henikoff,
164    2002) score since this was a lethal variant. SIFT predicts whether an amino acid substitution affects
165    protein function based on sequence homology and the physical properties of amino acids. The variant
166    impact categories are subjective agreements between VEPtools and SNPEff databases. However,
167    high-impact variants are considered to have protein level disruption or change, while modifier or
168    moderate variants impact non-coding regions of the genome. The SIFT score closer to zero is mostly
169    represented by HIGH or Modifier-impact categories, while tolerated levels (SIFT score of 0.05-1)
170    would show 'minimal' to 'no consequence' for the function of the genes under said variants. All sites
171    with high-impact scores were captured in a separate file for further processing.

172    **2.6    Novel variants**

173    The dbSNP database of NCBI (NCBI, 2019) contains data on variants already reported by
174    researchers. Novel variants are unlikely to be contained in these datasets and so failure to already
175    have a RS number, an indicator of being in a suitable database, may be another way to reduce the
176    search area along the genome.

177    A novel variant is unlikely to have been found already in previous studies so is likely to be found at
178    sites that are not already logged in the relevant SNP database. Sites in the final VCF file that did not
179    have an RS number were possible positions for the new variant. The VCF file was scanned for
180    positions not previously allocated a RS number using a script written in Perl 5.28. Such sites were
181    output for further analysis. The sites identified in this way were summarised by the number of
182    genotypes containing the variant allele found at each site. Candidate sites contained a 'potential'
183    variant allele in all nine samples.

184    **2.7    Comparing datasets**

185    Putting all these ideas together should make locating the site of the new variant on the genome
186    possible using WGS data from a small number of animals without the need for any other data
187    sources. The above analyses resulted in four independently-derived sets of data, each of which could
188    contain an indication of where the new variant might be found on the genome. These were 1) the sites
189    with the appropriate GCR, 2) sites in long ROH, with a high ABD score significant at $P < 0.01$, 3)
190    sites with a high-impact SIFT score and 4) sites with no RS number. Each dataset was derived from
191    the final VCF file by a method independent of the other three. If a site appears in all four datasets this
192    is likely to be the site of the new variant. The four datasets were compared for overlapping positions
193    by reading them into an Access database and linking on the site position.

194    **2.8    Predicting the effect of the novel variant**

195    The SIFT scoring method described above is one method for modelling the effects of a new variant
196    on the phenotype of the animal. PHYRE2 (Kelley et al., 2015) is an alternative approach which
197    searches for homologous sequences in a database of known proteins etc. The reference sequence and
198    its equivalent using the new variant were entered into the PHYRE2 database to see the effect of the
199    site of the proposed new variant on amino acid sequence and protein structure.

200    **2.9    Methods used to confirm the likely base position of the variant**

5

201 Two independent methods were used in an attempt to confirm the site derived from the methods used
202 on the WGS data described above. A sample of Irish Moiled animals comprising three bulls, 42 cows
203 and 18 dead calves (male and female) were genotyped using SNP-chips. These were then analysed
204 using 1) genotyping by PCR at the suggested site and 2) the ABD method on the SNP-chip data.

### 2.9.1 SNP processing

206 DNA samples extracted as described in the Supplementary File (Page S2) were genotyped in the
207 Department of Pathology, University of Cambridge (UK) and Gen-probe (Heron House, Oaks
208 Business Park, Crewe Road, Wythenshawe, Manchester, M23 9HZ, UK) using either: 1) the Illumina
209 BovineSNP50 BeadChip (Version 1, Illumina Inc., San Diego, CA) (50k SNP, n = 17); 2) the
210 BovineHD Genotyping BeadChip (777k SNP, n = 68) or 3) both chips (n = 7).

211 The SNP genotypes were prepared for all subsequent analyses using PLINK 1.9 (Chang et al., 2015;
212 PLINK, 2017; Purcell et al., 2007). Quality control parameters were used to edit the data. This
213 involved setting a lower limit on both sample and SNP quality at a call rate greater than 90%, and
214 SNPs were retained in the dataset if they were in Hardy-Weinberg Equilibrium. This was determined
215 using Fisher's Exact Test with a probability threshold of 0.05 and using the mid-p adjustment
216 described in (Graffelman and Moreno, 2013). The latest SNP positions were updated to the ARS-
217 UCD1.2 (GCF_002263795.1) build of the bovine genome using SNPchiMp (Nicolazzi et al., 2014;
218 Nicolazzi et al., 2015). In addition a merged set of data was produced using SNPchiMp combining all
219 genotyped animals from both LD and HD datasets with common SNP. This merged dataset was then
220 used in KING (Manichaikul et al., 2010) to generate relatedness coefficients between all genotyped
221 animals, based on whole genome SNP genotypes, to enable pedigree checking. Any animal whose
222 pedigree did not match the relatedness information from the SNP data was discarded. Because nine
223 animals were also used for the WGS analysis they too were excluded from the SNP ABD analysis, in
224 order to produce a dataset of independent animals.

### 2.9.2 Autozygosity by difference

226 The ABD method (Pollott, 2018), described above, was used on the merged SNP-chip dataset. The
227 probability of each SNP ABD score (difference between mean ROH length (Kb) from cases and
228 controls at each SNP position) was tested using 1,000 permutations of the dataset based on random
229 allocation of animals to phenotypes and recalculation of the ABD scores. Significance at the P <
230 0.001 level was considered as an indicator of a possible site of the new variant.

### 2.9.3 Genotyping by PCR analysis

232 Primers (5' CATGAACCCAGTGTCACAGC 3' and 5' CTCTCCGTGGAAGAGCAGAT 3') were
233 designed using Primer3 (version 4.1.0; http://primer3.ut.ee) to amplify a 218bp product spanning the
234 identified variant locus. Primer design was based on the published sequence for the *Bos taurus*
235 (UMD3.1; GCF_000003055.6) glucokinase gene ENSBTAG00000032288. Exon/intron boundaries
236 were derived from this in combination with mRNA RefSeq NM_001102302. PCR was performed
237 using AmpliTaq Gold polymerase (Applied Biosystems) according to the manufacturer's protocol.
238 Products were purified using the QIAquick PCR purification kit (Qiagen) and sequenced by Sanger
239 sequencing using the forward and reverse primers. Sequence analysis was carried out in CLC
240 Workbench. The candidate site was updated to the ARS-UCD1.2 (GCF_002263795.1) genome build
241 using the UCSC Genome LiftOver facility (UCSC, 2020).

### 3 Results

6

243  The WGS data from all nine animals, three cases and six controls, resulted in a final VCF file
244  comprising 8,234,367 biallelic autosomal single nucleotide variants which were to be used for all
245  subsequent WGS analyses. These are summarised by chromosome in Table 1 along with the length
246  of each chromosome aligned in the ARS-UCD1.2 (GCF_002263795.1) build of the cattle genome.

### 3.1    Genotype criteria

248  Searching the final VCF file for sites with the appropriate genotype criteria (all homozygous cases
249  for the same genotype and all controls heterozygous containing one allele forming the homozygote in
250  cases) resulted in 574 sites being identified. These are shown broken down by chromosome in Table
251  1. Applying the formula $1/3^n$ cases and $1/3^m$ parental controls to over 8 million SNV we would
252  expect to find ~418 sites fitting the genotype criteria. There were clearly more GCR sites than
253  expected in this set of animals. Chromosomes 1, 4, 12 and 23 appeared to have more sites than
254  expected (Table 1).

### 3.2    Autozygosity-by-difference method

256  In order to run the ABD method on the final VCF file the hard-filtering criteria shown in Table S1
257  were used on the extracted biallelic SNVs for the dataset. This resulted in a file of 629,716 SNP for
258  the ABD analysis. The ABD software was used to generate the Manhattan plots shown in Figures 2
259  and S4. The two plots in S4 show the mean ROH length at each of the base positions in the VCF file
260  for cases and controls respectively whilst Figure 2, the ABD score, shows the difference between
261  them. Long ROH were found on BTAs 4 and 18. Probabilities for the ABD scores were generated
262  from 100 permutations of the dataset and the regions of the genome with $P < 0.01$ are summarised in
263  Table S2. The 0.01 probability level was computed to be at an ABD score of 9,034 Kb. Table S2
264  shows that the length of BTA4 above the 0.01 probability threshold was 7.804Mb. The highest mean
265  ROH length in cases was 14.197Mb so the long ROH found on BTA4 continued either side of the
266  significant region. Similarly on BTA18, the highest mean ROH score in cases was 22.4Mb long.

### 3.3    SIFT score

268  Using the VEP to estimate the effect of each of the SNVs in the final VCF file resulted in 65,961
269  records of HIGH impact SNVs located at 7,764 different autosomal positions. The distribution of
270  these sites is summarised by chromosome in Table 1.

### 3.4    Sites with no RS number

272  The VCF file contained 340,893 sites with no RS number. Only 11% (6,298) sites had genotypes,
273  other than the homozygous reference genome, in all nine cases and controls (NoRS9). The
274  breakdown of these by chromosome is summarised in Table 1. These are likely to contain the novel
275  variant.

### 3.5    Overlap of GCR, ABD, NoRS9 and SIFT results

277  So far, four possible datasets were generated that might contain the site of the novel variant causing
278  this new autosomal recessive condition. The overlap between the four datasets is summarised in
279  Table 2. The genotype criteria method resulted in the fewest sites identified (574), with the other
280  methods increasing in the order autozygosity-by-difference, HIGH SIFT score and no RS number
281  with nine genotypes. Combining the GCR method with each of the others in turn allowed the
282  identification of 12 (ABD), 1 (HIGH SIFT) and 4 (NoRS9) sites in common. One site appeared in all
283  four datasets located at position Chr4: g.77173487A>T (ARS-UCD1.2 (GCF_002263795.1)). We
284  may tentatively conclude that this is the site of the novel variant causing early calf death in the Irish
285  Moiled breed.

286 **3.6 PHYRE2 pediction**

287 The PHYRE2 prediction (Kelley et al., 2015) of the secondary structures between the reference
288 genome and new variant *GCK* model at the beginning of Exon 8, which contains the possible site of
289 the new variant, Chr4: g.77173487A>T (ARS-UCD1.2 (GCF_002263795.1), predicted the amino
290 acid sequence NPGQQLWY from the reference genome being changed to NPGQQLLY with the new
291 variant.

292 **3.7 Independent confirmation of the results**

293 Over the period of work with Irish Moiled breeders a number of hair samples were collected and
294 these were genotyped on either: 1) the Illumina BovineSNP50 BeadChip (50k SNP, LD, n =17); 2)
295 the BovineHD Genotyping BeadChip (777k SNP, HD, n = 68) or 3) both chips (n = 7). This has
296 allowed us independently to assess the WGS results using two main approaches; Sanger sequencing
297 and the ABD runs of homozygosity method.

298 **3.7.1 Sanger sequencing**

299 A sample of 41 animals were Sanger sequenced at position Chr4: g.77173487A>T (ARS-UCD1.2
300 (GCF_002263795.1)) following PCR of the region surrounding this site. Table 3 shows the Fisher's
301 exact test (Fishre, 1922) results (with the Freeman-Halton extension (1951)) for animals falling into
302 three categories; calves, carriers and live animals of unknown status. The three genotypes are also
303 shown.

304 Table 3 shows that all TT animals were calves, all of which died of the symptoms described above.
305 All live animals were either AA or AT. The overall results were significant with a probability =
306 1.868e-05 (0.00001868) for this 3 x 3 table arising by chance, thus indicating a likely association
307 between genotype and health status at this site.

308 **3.7.2 ABD method based on genotypes derived from the PCR analysis**

309 The merged dataset (HD and LD chip data merged using SNPchiMp) comprised 63 animals (18 cases
310 and 45 controls) and 42,453 SNPs after quality control conditions were met. The 42 animals used in
311 the PCR analysis were selected for the ABD analysis, which excluded the WGS animals so that this
312 analysis was independent of the WGS ABD analysis. The animals were allocated to their case/control
313 status based on their PCR genotype at the highlighted location. The results of the ABD analyses are
314 summarised in Figure 3 and S5 and show a 20.8 Mb length of BTA4 with a permuted probability <
315 0.001, from 1,000 permutations, equivalent to an ABD score greater than 7,023 Kb. This region was
316 from position position Chr4: g. 62872037 to g. 83635054 (ARS-UCD1.2 (GCF_002263795.1))
317 which includes the site highlighted as the putative causal variant from the WGS analyses.

318 The results in Figure S5 show a long ROH on BTA21 but this was present in both the case and
319 control animals which negated each other in the ABD score analysis. This is a good example of the
320 benefit of the ABD method. Also, the long ROH found in WGS cases on BTA18 (Figure 2) was not a
321 feature of this larger set of results. There was reduced variability of these results with a higher
322 number of animals compared to those from the WGS dataset analysis with only nine animals.

323 **4 Discussion**

324 This work had two objectives; one general and the other more specific. The generally applied
325 objective was to test the idea that it is possible to find the site of a novel autosomal reccesive variant
326 using just a small number of whole-genome-sequenced animals and appropriate bioinformatic
327 methods, thus circumventing the need for the commonly-used two-stage approach highlighted by the

328  review of Pollott (2018) or the collection and/or use of further data. The specific objective was to
329  find the site of a new autosomal recessive condition thought to exist in Irish Moiled cattle.

330  **4.1    Bioinformatic methods used with WGS data to find the site of a new autosomal recessive**
331  **variant using a small number of cases and controls**

332  Whole genome sequencing is becoming more widely used to locate single novel variants with major
333  effects and a number of approaches have been used. In a large scale analysis of over Holstein cattle
334  WGSs seven dominant conditions were located using the genome criteria approach (Bourneuf et al,
335  2017) involving one case for each of the seven conditions and a control population of 1,230 animals.
336  The trio approach has been used by a number of authors (see for example Sayyab et al., 2016). This
337  method takes WGSs of one affected offspring and its two parents and uses the genotype criteria
338  method to find possible sites for the causative variant. The large number of sites identified are further
339  reduced by a range of methods. Using one dog example (Sayyab et al., 2016), a filtering pipeline was
340  established with 7 steps, including genotype criteria and SIFT analysis, Sanger sequencing
341  verification and sequencing of an additional 24 cases/controls. Runs of homozygosity methods have
342  been widely used with SNP-chip data (Pollott, 2018) and Letko et al. (2020) report an example of
343  using this method in Zwartbles sheep to locate a novel autosomal recessive condition associated with
344  Type 1 Primary Hyperoxaluria. Their study relied upon additional data from both the Sheep
345  Genomes project and 79 publicly available genomes of various breeds to provide 'control' data for
346  the GCR method.

347  In the current study four bioinformatic methods were tested to try to find the location of the new
348  variant causing early calf death in the nine Irish Moiled animals and relied on the 'correct' site
349  appearing in all four methods. No additional data from the Irish Moiled or any other breeds were
350  used. Two of the methods (GCR and ABD) do not require any prior information about genes, SNP or
351  other genomic features but rely on across- and within-animal patterns of information contained in the
352  genotypes at each SNV found in the VCF file. Ideally one would like to use these two methods alone
353  since they are not only independent of any prior knowledge about genome features (except the
354  reference genome for the alignments and generation of the VCF file) but they will also be able to find
355  a new variant causing an autosomal recessive condition anywhere on the genome, even when located
356  outsde a protein coding region: a useful feature of the two methods. As has been seen, using three
357  cases and six controls with the ABD method and GCR combined revealed 12 possible sites in a 7.795
358  Gb stretch of Chr4 between g.70889821 and g.78684588 (ARS-UCD1.2 (GCF_002263795.1))
359  involving 18,705 SNVs. The underlying implication of this approach is that with more animals,
360  either cases or controls, we would find fewer sites and so make the search somewhat more
361  straightforward and find just a single causative variant site.

362  **4.1.1   The genotype criteria approach**

363  The method used to find the sites meeting the genotype criteria was based on a number of implied
364  assumptions not stated by Pollott (2018). Firstly that GCR sites would be evenly distribution across
365  the genome. Secondly, the higher the number of animals used the greater the chances of finding the
366  GCR site of the new variant. Thirdly, a GCR site was not dependent on the balance of cases and
367  controls in the samples. Fourthly, the location of a single GCR site was independent of the genetic
368  relationship between cases and controls. Each assumption was tested using the data analysed in this
369  study, either the final VCF file for BTA4 or the SNP-chip data with phenotypes allocted by the PCR

370    results as appropriate. The detail of these investigations is given in the Supplementary File (Pages
371    S10-17).

**4.1.1.1 Evenly-spaced GCR sites across the genome**

373    The results in Table 1 show that some chromosomes contain no GCR candidate sites at all (BTAs 3,
374    7, 8, 9, 10, 13, 15, 16, 18, 19, 20, 21, 22, 25, 27, 28 and 29). In fact these 17 chromosomes comprise
375    55% of the autosomal genome. Many chromosomes contained far fewer GCR sites than expected
376    whereas others contained a much greater number than expected (BTAs 1 and 23). A GCR site (in this
377    case) comprises two components; the 0/1 in all controls and the 1/1 in all cases (using 0 to mean the
378    reference allele and 1 the new or alternative variant). We might expect the chromosomes containing
379    long ROH in cases potentially to have many more GCR sites than others. Inspection of Figure S4
380    shows BTAs 4 and 18 as having the longest mean ROH but only BTA1 had a large excess of GCR
381    sites. Using the information shown in the Supplementary File, the original implied assumption from
382    Pollott (2018) of an even distribution of GCR sites across the genome has not been verified here and
383    this has implications for the number of animals required to find a new variant site using WGS data
384    alone. Clearly the long ROH on BTA 4 were linked to a large number of GCR sites but that on
385    BTA18 was not.

**4.1.1.2 Using more animals increases the chances of finding a causative  GCR site**

387    The second assumption about the use of GCR sites to locate a novel autosomal recessive variant was
388    that the greater the number of animals that are genotyped the better the chance of locating the new
389    variant. It appears, from the work reported in the Supplementary File, that the number of genotyped
390    animals required to find the single candidate sites is 3 to 4 times greater than predicted using the
391    original formula of Pollott (2018); one case and 29 controls appears to require the fewest total
392    animals genotyped to find the single candidate site in the SNP-chip dataset used. However, this 'long
393    tail' is due to several GCR sites being close together around the candidate site and always being
394    'found' in the generated datasets. Differentiating between them may require another method or using
395    a different set of controls, perhaps from more distantly related individuals.

**4.1.1.3 The balance of cases and controls**

397    The basic calculation of the number of animals required to find a GCR site is independent of the
398    balance between the number of cases and controls used. Figure S10 demonstrates that the lower the
399    number of cases used the fewer the total number of animals required to be genotyped. At first sight
400    these are rather startling results. However, outside the candidate site, it is much more unlikely to find
401    all controls with a heterozygote genotype, whereas there will be many sites with all homozygous
402    genotypes in cases; after all long ROH imply many 1/1 genotypes and so more sites potentially could
403    meet the genotype criteria. The information in the Supplementary File clearly demonstrates that the
404    number of animals required is much closer to the theoretical numbers when large ROH regions are
405    excluded from the analyses. The number of animals required to find the candidate site is still slightly
406    greater than the theoretical figure but this may be due to some other small areas of ROH not removed
407    from the dataset. There are an enduring number of GCR sites in the high-ROH regions which inflate
408    the results in contrast to the theoretical number of sites expected. This illustrates why theory and
409    'practice' may differ.

**4.1.1.4 The genetic relationship between cases and controls**

10

411　The nine animals used in the WGS analysis comprised three cases and six parental controls. In order
412　to investigate whether the closely-related controls might inflate the number of GCR sites found, an
413　alternative SNP-chip dataset was derived using controls that were most distantly related to the cases
414　(Supplementary File and Figure S12). In this case the number of animals required to find the new
415　variant site was much closer to the theoretical expectation than with the parental controls.

### 4.1.1.5 Final genotype criteria summary

417　In summary, the method of locating the site of a new autosomal recessive variant using genotype
418　criteria has potential. The theoretical expectation of the number of animals needed to be genotyped
419　underestimates the actual number required if the method was the sole way to find the new variant.
420　The 'simulated' datasets in the Supplementary File indicated that the outcome of any *single* set of
421　genotyped animals is likely to vary widely, so needing fewer animals is just as likely as needing
422　more, but of course it is unpredictable for any given set of animals. The best chance of success is
423　likely to come from genotyping slightly more than the theoretical number of animals, using 33-50%
424　from cases and the remainder from more distantly related controls (i.e. any animals from the breed).
425　Relating this back to the 574 GCR sites found from the three cases and six parental controls used in
426　the WGS analyses, here the number of cases used was at the lower bound of this recommendation
427　(33%). The excess of GCR sites found was likely to have been caused by using parental controls,
428　although it was well within the bounds of variability estimated from the SNP data.

### 4.1.2　The Autozygosity-By-Difference method

430　The ABD method was developed originally to locate regions of the genome likely to contain a new
431　autosomal recessive variant using SNP data (see for example Posbergh et al., 2018). In the current
432　analysis it has been applied to WGS data for the first time, as well as being used to confirm the
433　results in an independent sample of SNP-genotyped animals. Because the method is sensitive to
434　incorrectly called genotypes, a feature of WGS data, it was necessary to employ hard filtering criteria
435　(see Table S1) of both sites and genotypes in order to get a useable set of data. In this case the VCF
436　file was reduced from ~12 million to ~630,000 sites but this would differ under alternative hard-
437　filtering criteria. Since VQSR methods were not available in the current situation (due to the species
438　and number of animals genotyped) an alternative approach was taken; selecting sites and genotypes
439　to fall with ±2 s.d. of the mean (or peaks in the case of bimodal variables). This allowed the location
440　of several long ROH, one of which was found, by additional methods, to contain the new variant.

441　As well as using alternative hard-filtering criteria it may be possible to use other approaches,
442　including site sampling or sliding-windows, to locate the region containing the new autosomal
443　recessive variant. Using a site-sampling approach with a VCF file one could randomly select, say, 5-
444　10% of sites evenly spread across the genome with the ABD method. Repeated samples of these
445　SNVs, (say 100), could be randomly drawn and the 100 sets of ABD results averaged at each site.
446　Alternatively, one could use a sliding window of, say, 10,000 base positions and count the number of
447　homozygous variant case and control genotypes in each window. The window would then be moved
448　along the genome at a given interval, say every 1,000 base positions, and the results plotted. One
449　would expect to find the new variant causing the autosomal recessive condition in the region with the
450　highest ABD-type score. Both these methods would overcome the problem of a single incorrectly-
451　called genotype disrupting the long ROH in the ABD results and the need for hard filtering.

452 Using ABD on the hard-filtered WGS data resulted in the identification of two regions of the genome
453 having an ABD score above the 0.01 probability threshold, and therefore likely to contain the new
454 variant (Figure 2). Two aspects of Figures 2 and S4 are of note. Firstly there were a number of long
455 ROH found in the controls throughout the genome with BTA18 having the largest mean ROH length.
456 This was also found in the cases but the effect of combining the two sets of data in the ABD score
457 was to remove many of these 'breed-specific' ROH and leave those which probably harboured the
458 new variant. This is one of the advantages of the ABD method, particularly in rare breeds, but it has
459 also been shown to be effective in removing long ROH associated with new variants in the myostatin
460 gene in both Texel sheep (Pollott, 2013) and Piedmontese cattle (Biscarini et al., 2013).

461 The ABD method was also used for another purpose in this work; to confirm the WGS results on an
462 independent set of animals with SNP-chip-derived data (Figure 3). As with the WGS ABD results,
463 there were many long ROH in the SNP-chip dataset in both cases and controls. In this instance there
464 were two very long ROH on BTA4 and BTA21, but interestingly not BTA18 as was found in the
465 WGS data. The smaller number of animals used in the WGS analysis probably resulted in BTA18
466 having a long ROH due to sampling of closely related individuals. Figure S5 also shows BTA18 to
467 have a long ROH but it was less pronounced in this larger dataset. In Figure 3 the result of
468 subtracting the control ROHs from that of cases at each site reduced the noise considerably and left
469 BTA4 as the only significant peak by a considerable margin. Once again it has been demonstrated
470 that the power of the ABD method to remove noise works effectively.

471 **4.1.3 Combining genotype criteria and ABD results**

472 The approach in this work has been to use several bioinformatic methods on WGS data to see if they
473 can pinpoint the site of a causal variant of a new autosomal recessive condition. Table 2 has
474 highlighted a significant region on BTA4 using the ABD method that contained 12 sites meeting the
475 genotype criteria. The discussion above has suggested that using more animals may have reduced the
476 number of candidate sites by a small amount but a greater use of unrelated controls may have reduced
477 the number of GCR sites in the target area more effectively.

478 **4.1.4 SIFT score**

479 In this set of results the SIFT scores were the crucial factor in determining the site of the new variant.
480 Position Chr4: g.77173487A>T (ARS-UCD1.2 (GCF_002263795.1) was the only one of the 18 GCR
481 sites in the target area to have a high-impact SIFT score (0-0.05). However, only sites in or near a
482 coding region are scored using the SIFT method so it is not always going to find the causative site if
483 it is located outside these regions of the genome.

484 **4.1.5 RS number**

485 The use of the absence of a RS number could be useful but, in this case, did not prove to be the final
486 factor locating the novel variant site.

487 **4.2 Other types of inheritance and effects**

488 This paper reports the search for a new autosomal recessive variant causing a fatal condition in calves
489 using a range of bioinformatic methods. It raises two issues relating to non-fatal autosomal recessive
490 conditions and to other modes of inheritance.

491 **4.2.1 Non-fatal conditions**

12

492    There is nothing particularly special about the search for the causal variant of new fatal genetic
493    diseases compared to non-fatal conditions. Clearly the potential phenotype is obvious but it may be
494    harder in the initial stages of a new disease to correctly phenotype the dead animals. With a non-fatal
495    condition it is still necessary to correctly phenotype the animals but there are likely to be more
496    opportunities to test them and to return to them for samples, in many instances.

497    **4.2.2   Other modes of inheritance**

498    Having suggested this WGS approach for finding a new autosomal recessive variant the question
499    arises about its usefulness with variants involving other modes of inheritance.

500    A dominance mode of inheritance can be thought of as the reverse of the recessive mode. One would
501    expect to find cases to be 0/1 or 1/1 and controls to be 0/0 so the genotype criteria would be different
502    compared to the recessive case. However, the number of animals required to use the GCR method
503    would be very similar with cases being $2/(3^n)$ and controls $1/(3^m)$; the numerator having little effect
504    with such a large denominator. The ABD method could only be used if all cases were 1/1 but that is
505    unlikely with a dominant condition due to the large number of heterozygotes likely to be in the
506    population. Alternatively, if there was some way to phenotypically distinguish 1/1 from 1/0 cases this
507    would be useful. The 0/0 controls are unlikely to be situated in long ROH since they are likely to
508    have been subjected to many generations of recombination, so alternative methods may be required.
509    The new dominant variant would be situated in a long haplotype so it may be possible to adapt
510    haplotype discovery methods to this situation. Both the SIFT score and RS number methods would
511    be applicable but they are less powerful than the other two both because they rely on previous
512    knowledge and, in the case of SIFT, it only works for a limited distance around a protein-coding
513    region.

514    These methods could be used for a recessive sex-linked new variant i.e. one found on the X
515    chromosome. Males would provide no useful data in this case so only females would be required.
516    Both the ABD and the GCR methods would work the same way but with a lot fewer sites to search
517    (only the X chromosome data would be needed).

518    **4.3    Finding the causative variant for a perinatal mortality syndrome in Irish Moiled Cattle**

519    The likely site for the causative variant of this fatal perinatal condition in Irish Moiled animals has
520    been successfully located using just six parental control animals and three cases. Perinatal mortality
521    (within 24 h of birth) typically occurs in about 6-10% of calves born (Brickell et al., 2009) with a
522    further 3-4% dying in their first month, mainly from infectious disease (Johnson et al., 2017). The
523    site highlighted at Chr4: g.77173487A>T (ARS-UCD1.2 (GCF_002263795.1) was located in the
524    glucokinase gene (*GCK*) and is a splice acceptor variant. Analysis of the OMIA website (OMIA,
525    2020; Nicholas and Hobbs, 2012) showed splice acceptor variants to be responsible for ~8% of
526    known variants in non-laboratory animals. There was a clear difference in the PHYRE2 prediction of
527    the secondary structures between the reference genome and new variant *GCK* model. As observed
528    from the SIFT-score results this is expected to have a disruptive effect on the operation of the *GCK*
529    gene.

530    Glukokinase is a key enzyme found in the liver, pancreas, brain and endocrine cells of the gut. It
531    catalyses the starting point of glycolysis by phosphorylating glucose to form glucose-6-phosphate
532    (Matschinsky et al., 1993). The crystal structure has revealed that glucose binds in a deep cleft

533    between a large and small domain of GCK, resulting in a conformational change and enzyme
534    activation (Kamata et al., 2004). Glucokinase stimulates glucose uptake, glycolysis and glycogen
535    synthesis by hepatocytes, whereas in pancreatic β-cells it plays a crucial role in glucose-stimulated
536    insulin secretion. Glucose homeostasis is essential in mammals and is under tight endocrine control,
537    with insulin acting as the key regulator.

538    There are currently  922 SNPs listed within the bovine *GCK* gene (NCBI, 2019) but the closest to the
539    new variant site flanking either side were at Chr4: g.77173441 (an intron variant) and Chr4:
540    g.77174392 (ARS-UCD1.2 (GCF_002263795.1)), some 46 and 905 bp away respectively. A segment
541    of the ARS-UCD1.2 (GCF_002263795.1) genome 30 bp either side of the candidate variant was
542    selected and blastn (Altschul et al., 1990) was used with the 61 bp sequence to find any homologous
543    region on the human genome. A 40 bp length of sequence was found with 35 identical bases and a
544    score of 50.9 bits (55) and no gaps. This was located on the reverse strand of human Chr7:
545    g.44146590 to g.44146629 (GRCh38.p13 (GCF_000001405.39), in the *GCK* gene. The location on
546    the human genome equivalent to the candidate variant found in Irish Moiled calves was at Chr7:
547    g.44146620 (GRCh38.p13 (GCF_000001405.39). This was a highly conserved site with 96 out of the
548    100 vertebrate genomes shown on the UCSC (UCSC, 2020) genome browser all having a T on the
549    forward strand, the remaining 4 being not reported. No SNP was found at this site in the human
550    database but there was a SNP reported at the adjacent position (Chr7: g.44146619; GRCh38.p13
551    (GCF_000001405.39)) which was catalogued as rs1167675604, a C>T change on the forward strand.
552    This site was also highly conserved in 96 out of the 100 vertebrate genomes on the UCSC genome
553    browser and was also a splice-site acceptor variant. The ClinVar (ClinVar, 2020) record for this
554    variant states that "The variant disrupts a canonical splice site, and is therefore predicted to result in
555    the loss of a functional protein. Found in at least one symptomatic patient, and not found in general
556    population data." It's incidence was estimated to be well below 0.001% of the population. In addition
557    the Varsome (Varsome, 2020) record for this SNP states that the effect of the variant was 'Very
558    Strong' which means " Null variant (intronic within ±2 of splice site) affecting gene *GCK*, which is a
559    known mechanism of disease (gene has 378 known pathogenic variants which is greater than
560    minimum of 3), associated with Diabetes mellitus, permanent neonatal 1, Maturity-onset diabetes of
561    the young, type 2 and Hyperinsulinemic hypoglycemia, familial 3."

562    The mouse genome was also investigated in the same way but no SNP were found in the candidate
563    region.

564    Over 600 variants have been reported in the human *GCK* gene, which have varying effects depending
565    on their location (Osbak et al., 2009; OMIM 138079). Heterozygous inactivating variants cause a
566    condition known as maturity onset diabetes of the young, characterised by mild fasting
567    hyperglycaemia. Homozygotes are much more rare in the human population, and neonates present
568    earlier with permanent neonatal diabetes mellitus. In mice, however, pups born with global *GCK*
569    knockout (-/-) are slightly smaller than wildtype animals (+/+), have glucose levels about eightfold
570    higher and die within 3 to 5 days (Grupe et al., 1995). Tissue specific β-cell knockouts die within 4
571    days of birth whereas hepatic knockout impairs glucose utilization and glycogen synthesis but with
572    only mild hyperglycaemia (Postic et al., 1999).

573    Pregnancy outcome in women depends on a combination of the genotype of both mother and fetus
574    (Spyer et al., 2001). When the fetus carries a single *GCK* variant this affects glucose homeostasis
575    with reduced insulin secretion, so both placental and birth weight are reduced (Hattersley et al., 1998;
576    Spyer et al., 2008). During pregnancy, the fetal glucose supply is derived almost entirely from the

577    dam across the placenta using facilitated diffusion by glucose transporters. In ruminants this uptake is
578    regulated sequentially by GLUT1 and GLUT3 (SLC2A1 and SLC2A3) (Wooding et al., 2005).

579    The fetus has a low capacity for endogenous glucose production but this increases in late gestation, in
580    response to the pre-term rise in glucocorticoid production, together with catecholamine and thyroid
581    hormone stimulation. These promote hepatic glycogen synthesis and gluconeogenesis, which are
582    essential in providing the neonatal calf with an adequate glucose supply as milk lactose on its own is
583    insufficient (Hammon et al., 2013). The postnatal maturation in the regulation of energy supply may
584    thus explain why lack of GCK activity is fatal at this stage of life.

585    **4.4    Concluding remarks**

586    The original intention for this work was to locate the site of a novel variant causing perinatal
587    mortality in Irish Moiled calves. This has been achieved, and shown to be located in the *GCK* gene,
588    but in the process it became apparent that there were no straightforward ways to achieve this
589    objective. At best, a two-stage approach was required involving genotyping a group of cases and
590    controls, identifying the genomic region likely to contain the novel variant followed by further work
591    to sequence the identified region and look for appropriate signals in the data. Consequently, a further
592    objective was set in order to simplify the process and investigate whether it would be possible to use
593    a single whole genome sequencing stage with appropriate bioinformatic methodology to find the
594    candidate site. This, too has been achieved by sequencing nine animals, three cases and six parental
595    controls, and applying four methods to the data. In the process it has been possible to investigate
596    some of these methods in more detail and arrive at some general conclusions to aid future such
597    studies.

598    The VCF file format has proved to be a very practical source of data for this study particularly
599    because it reduced the search 'area' from over 2.5 billion base positions down to one involving only
600    8 million sites. In addition, the VCF file format facilitated finding the novel site when combined with
601    methods to interrogate it for genotype criteria, long runs of homozygosity and the predicted effects of
602    variants on the phenotype of the animal. Using these three methods allowed the identification of a
603    single variant site which was found to have both the genomic and biological properties associated
604    with this novel condition.

605    In the process of carrying out this work it has been possible to refine the genotype criteria method to
606    demonstrate that in reality only a small number of cases and controls *is* required, controls should
607    outnumber cases by 2:1 and controls should be more distantly related to cases. In addition it has been
608    possible to show that using a runs-of-homozygosity method, previously only used on SNP-chip
609    genotype data, with whole-genome-sequence data it was possible to locate the region of the genome
610    containing the novel variant.

611    In future it should be possible to use the combination of genotype criteria and runs-of-homozygosity
612    methods with the appropriate number of cases and controls, suitably distantly related, to locate the
613    site of any new autosomal recessive genetic condition in a relative short time. This should then
614    facilitate a more speedy elimination of the harmful variant from the population by using an
615    appropriate genetic test on available animals.

616

617    **Conflict of Interest**

620 **Author Contributions**

621 GP and DCW designed the study. GP analysed the data, wrote the ABD software, carried out the
622 bioinformatic analyses and wrote the first draft of the paper. MS produced the alignments, VCF files,
623 SIFT and PHYRE2 scores. RP and CM carried out the Sanger sequencing. GP, DCW, RP and MS
624 contributed written material to the final paper.

636 **Data Availability Statement**

637 The data is owned by the Irish Moiled Breeders.

638 **References**

639 Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. (1990). Basic local alignment
640 search tool. J. Mol. Biol. 215, 403-410. doi: 10.1016/S0022-2836(05)80360-2

641 Biscarini, F., Corvo, M., del, Stella, A., Albera, A., Ferenčakovic, M., and Pollott, G.E. (2013).
642 Looking for the mutations for arthrogryposis and macroglossia in Piedmontese cattle:
643 preliminary results. J. sobre Prod. Anim., Zaragoza 14 y 15 de mayo de 2013. 538-540.

644 Bourneuf, E., Otz, P., Pausch, H., et al. (2017). Rapid discovery of de novo deleterious mutations in
645 cattle enhances the value of livestock as model species. Sci. Rep. 7, 11466. doi:
646 10.1038/s41598-017-11523-3.

647 Brickell, J.S., McGowan, M.M., Pfeiffer, D.U., and Wathes, D.C. (2009). Mortality in Holstein-
648 Friesian calves and replacement heifers, in relation to body weight and IGF-I concentration, on
649 19 farms in England. Anim. 3, 1175-82. doi: 10.1017/S175173110900456X.

650    Broad Institute (2020). https://gatk.broadinstitute.org/hc/en-us. [Accessed July 26, 2021].

651    Chang, C.C., Chow, C.C., Tellier, L.C.A.M., Vattikuti, S., Purcell, S.M., and Lee, J.J. (2015).
652        Second-generation PLINK: rising to the challenge of larger and richer datasets. GigaScience. 4,
653        s13742-015-0047-8. doi: https://doi.org/10.1186/s13742-015-0047-8.

654    ClinVar. (2020). https://clinvarminer.genetics.utah.edu/ [Accessed October 20, 2020].

655    Danecek, P., Auton, A., Abecasis, G., Albers, C.A., Banks, E., DePristo, M.A., Handsaker, R.E.,
656        Lunter, G., Marth, G.T., Sherry, S.T., McVean, G., and Durbin, R. (2011). The variant call
657        format and VCFtools. Bioinformatics (Oxford, England). 27, 2156-8. doi:
658        10.1093/bioinformatics/btr330. Epub 2011 Jun 7.

659    Ensembl. (2020) http://oct2018.archive.ensembl.org/Bos_taurus/Info/Annotation [Accessed October
660        20, 2020].

661    Fisher, R. A. (1922). On the interpretation of χ2 from contingency tables, and the calculation of P. J.
662        R. Stat. Soc. 85, 87–94. doi:10.2307/2340521. JSTOR 2340521.

663    Freeman, G.H., and Halton, J.H. (1951). Note on exact treatment of contingency, goodness of fit and
664        other problems of significance. Biometrika. 38, 141-149.

665    Graffelman, J., and Moreno, V. (2013). The mid p-value in exact tests for Hardy-Weinberg
666        equilibrium. Stat. Appl. Genet. Mol. Biol. 12, 433-448. doi:10.1515/sagmb-2012-0039

667    Grupe, A., Hultgren, B., Ryan, A., Ma, Y.H., Bauer, M., and Stewart, T,A. (1995), Transgenic
668        knockouts reveal a critical requirement for pancreatic beta cell glucokinase in maintaining
669        glucose homeostasis. Cell. 83, 69-78. doi: 10.1016/0092-8674(95)90235-x.

670    Hammon, H.M., Steinhoff-Wagner, J., Flor, J., Schönhusen, U., and Metges, C.C. (2013). Lactation
671        Biology Symposium: role of colostrum and colostrum components on glucose metabolism in
672        neonatal calves. J. Anim. Sci. 91, 685-95. doi: 10.2527/jas.2012-5758.

673    Hattersley, A.T., Beards, F., Ballantyne, E., Appleton, M., Harvey, R., and Ellard, S. (1998).
674        Mutations in the glucokinase gene of the fetus result in reduced birth weight. Nat. Genet. 19,
675        268-70. doi: 10.1038/953.

676    Irish Moiled Cattle Society. (2020). https://www.irishmoiledcattlesociety.com/breed-history/.
677        [Accessed October 20, 2020]

678    Johnson, K.F., Chancellor, N., Burn, C.C., and Wathes, D.C. (2017). Prospective cohort study to
679        assess rates of contagious disease in pre-weaned UK dairy heifers: management practices,
680        passive transfer of immunity and associated calf health. Vet. Rec. Open. 4, e000226. doi:
681        10.1136/vetreco-2017-000226.

682    Kamata, K., Mitsuya, M., Nishimura, T., Eiki, J., and Nagata, Y. (2004). Structural basis for
683        allosteric regulation of the monomeric allosteric enzyme human glucokinase. Structure. 12,
684        429-38. doi: 10.1016/j.str.2004.02.005.

685    Kelley, L., Mezulis, S., Yates, C. et al. (2015). The Phyre2 web portal for protein modeling,
686        prediction and analysis. Nat Protoc. 10, 845–858. https://doi.org/10.1038/nprot.2015.053

687    Letko, A., Dijkman, R., Strugnell, B., Häfliger, I.M., Paris, J.M., Henderson, K., Geraghty, T., Orr,
688        H., Scholes, S., and Drögemüller, C. (2020). Deleterious AGXT missense variant associated
689        with Type 1 primary hyperoxaluria (PH1) in Zwartbles sheep. Genes (Basel). 11, 1147.
690        https://doi.org/10.3390/genes11101147.

691    Manichaikul, A., Mychaleckyj, J.C., Rich, S.S., Daly, K., Sale, M., and Chen, W.M. (2010). Robust
692        relationship inference in genome-wide association studies. Bioinformatics. 26, 2867-2873. doi:
693        10.1093/bioinformatics/btq559.

694    Matschinsky, F., Liang, Y., Kesavan, P., Wang, L., Froguel, P., Velho, G., Cohen, D., Permutt, M.A.,
695        Tanizawa, Y., Jetton, T.L., *et al*. (1993). Glucokinase as pancreatic beta cell glucose sensor and
696        diabetes gene. J. Clin. Invest. 92, 2092–2098. doi: 10.1172/JCI116809.

697    McLaren, W., Gil, L., Hunt, S.E., Riat, H.S., Ritchie, G.R., Thormann, A., Flicek, P., and
698        Cunningham, F. (2016). The Ensembl Variant Effect Predictor. Genome Biol. 17, 122.
699        doi:10.1186/s13059-016-0974-4.

700    NCBI (2019). https://www.ncbi.nlm.nih.gov/snp/. [Accessed October 20, 2020]

701    Ng, P.C., and Henikoff, S. (2002). Accounting for human polymorphisms predicted to affect protein
702        function. Genome Res. 12: 436-446. doi: 10.1101/gr.212802.

703    Nicholas, F.W., and Hobbs, M. (2013). Mutation discovery for Mendelian traits in non-laboratory
704        animals: a review of achievements up to 2012. Anim. Genet. 45,157-70. doi:
705        10.1111/age.12103.

706    Nicolazzi, E.L., Caprera, A., Nazzicari, N., Cozzi, P., Strozzi, F., Lawley, C., Pirani, A., Soans, C.,

707        Brew, F., Jorjani, H., Evans, G., Simpson, B., Tosser-Klopp, G., Brauning, R., Williams, J.L.,

708        and Stella A. (2015). SNPchiMp v.3: integrating and standardizing single nucleotide

709        polymorphism data for livestock species. BMC Genomics, 16:283. doi:

710        https://doi.org/10.1186/s12864-015-1497-1

711    Nicolazzi, E.L., Picciolini, M., Strozzi, F., Schnabel, R.D., Lawley, C., Pirani, A., Brew, F. and Stella

712        A. (2014). SNPchiMp: A database to disentangle the SNPchip jungle in bovine livestock. BMC

713        Genomics. 15, 123. doi: 10.1186/1471-2164-15-123.

714    OMIA. (2020). Online Mendelian Inheritance in Animals, OMIA. Sydney School of Veterinary

715        Science, https://omia.org/. [Accessed October 20, 2020].

716    Osbak, K.K., Colclough, K., Saint-Martin, C., Beer, N.L., Bellanné-Chantelot, C., Ellard, S., and

717        Gloyn, A,L. (2009), Update on mutations in glucokinase (GCK), which cause maturity-onset

718        diabetes of the young, permanent neonatal diabetes, and hyperinsulinemic hypoglycemia. Hum.

719        Mutat. *3*0, 1512-26. doi: 10.1002/humu.21110.

720    PLINK. (2017) https://zzz.bwh.harvard.edu/plink/contact.shtml. [Accessed 20 October, 2020].

721    Pollott, G.E. (2013). Do selective sweeps in sheep breeds indicate the genomic sites of breed

722        characteristics? Book of abstracts of the 64th European Association for Animal Production

723        Annual Meeting, Nantes, France. p627. Wageningen Academic Publishers; Netherlands.

724    Pollott, G. E. (2018). Invited review: Bioinformatic methods to discover the likely causal variant of a

725        new autosomal recessive genetic condition using genome-wide data. Anim. 12, 2221–2234.

726        doi: 10.1017/S1751731118001970.

727    Posbergh, C.J., Pollott, G.E., Southard, T.L., Divers, T.J., and Brooks, S.A. (2018). A non-

728        synonymous change in adhesion G protein-coupled receptor L3 associated with risk for Equine

729        Degenerative Myeloencephalopathy in the Caspian Horse. J. Equine. Vet. Sci. 70, 96-100.

730    Postic, C., Shiota, M., Niswender, K.D., Jetton, T.L., Chen, Y., Moates, J.M., Shelton, K.D., Lindner,

731        J., Cherrington, A.D., and Magnuson, M.A. (1999). Dual roles for glucokinase in glucose

732        homeostasis as determined by liver and pancreatic beta cell-specific gene knock-outs using Cre

733        recombinase. J. Biol. Chem. 274, 305-15. doi: 10.1074/jbc.274.1.305.

734 Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M.A.R., Bender, D., Maller, J., Sklar,
735    P., de Bakker, P.I.W., Daly, M.J., and Sham, P.C. (2007). PLINK: a toolset for whole-genome
736    association and population-based linkage analysis. Am. J. Hum. Genet. 81, 559-75. doi:
737    10.1086/519795.

738 Sayyab, S., Viluma, A., Bergvall, K., Brunberg, E., Jagannathan, V., Leeb, T., Andersson, G., and
739    Bergström, T.F. (2016). Whole-Genome sequencing of a canine family trio reveals a FAM83G
740    variant associated with hereditary footpad hyperkeratosis. G3 (Bethesda). 6, 521-7. doi:
741    10.1534/g3.115.025643.

742 Spyer, G., Hattersley, A.T., Sykes, J.E., Sturley, R.H., MacLeod, and K.M. (2001). Influence of
743    maternal and fetal glucokinase mutations in gestational diabetes. Am. J. Obstet. Gynecol. 185,
744    240-1. doi: 10.1067/mob.2001.113127.

745 Spyer, G., Slingerland, A.S., Knight, B.A., Ellard, S., Clark, P.M., Hauguel-de Mouzon, S., and
746    Hattersley, A,T. (2008), Mutations in the glucokinase gene of the fetus result in reduced
747    placental weight. Diabetes Care. 31, 753-7. doi: 10.2337/dc07-1750.

748 UCSC. (2020) https://genome.ucsc.edu/. [Accessed October 20, 2020].

749 Varsome. (2020) https://varsome.com/variant/hg19/rs1167675604. [Accessed October 20, 2020].

750 Wooding, F.B., Fowden, A.L., Bell, A.W., Ehrhardt, R.A., Limesand, S.W., and Hay, W.W. (2005).
751    Localisation of glucose transport in the ruminant placenta: implications for sequential use of
752    transporter isoforms. Placenta. 26, 626-40. doi: 10.1016/j.placenta.2004.09.013.

753

754

755 **Table 1** **A summary of results by chromosome**

756

| BTA | Length (bp) | Number of SNV in VCF file | Estimated number of GCR | Actual number of GCR | Number of SIFT sites | Number of NoRS9 sites |
|-----|-------------|---------------------------|-------------------------|----------------------|----------------------|-----------------------|
| 1 | 158,534,110 | 535,135 | 27 | 489 | 251 | 648 |
| 2 | 136,231,102 | 452,534 | 23 | 2 | 279 | 164 |
| 3 | 121,005,158 | 394,128 | 20 | 0 | 383 | 271 |
| 4 | 120,000,601 | 373,759 | 19 | 13 | 485 | 342 |
| 5 | 120,089,316 | 401,056 | 20 | 1 | 458 | 173 |
| 6 | 117,806,340 | 362,930 | 18 | 2 | 174 | 220 |
| 7 | 110,682,743 | 358,947 | 18 | 0 | 491 | 249 |
| 8 | 113,319,770 | 319,261 | 16 | 0 | 204 | 485 |
| 9 | 105,454,467 | 328,559 | 17 | 0 | 179 | 184 |
| 10 | 103,308,737 | 353,791 | 18 | 0 | 293 | 161 |
| 11 | 106,982,474 | 324,386 | 16 | 6 | 337 | 132 |
| 12 | 87,216,183 | 335,302 | 17 | 15 | 139 | 145 |
| 13 | 83,472,345 | 232,616 | 12 | 0 | 249 | 188 |
| 14 | 82,403,003 | 262,990 | 13 | 1 | 108 | 135 |
| 15 | 85,007,780 | 285,748 | 15 | 0 | 374 | 158 |
| 16 | 81,013,979 | 268,622 | 14 | 0 | 225 | 388 |
| 17 | 73,167,244 | 278,406 | 14 | 1 | 198 | 106 |
| 18 | 65,820,629 | 194,623 | 10 | 0 | 522 | 384 |
| 19 | 63,449,741 | 227,124 | 12 | 0 | 449 | 110 |
| 20 | 71,974,595 | 230,092 | 12 | 0 | 110 | 84 |
| 21 | 69,862,954 | 202,249 | 10 | 0 | 205 | 260 |
| 22 | 60,773,035 | 180,922 | 9 | 0 | 147 | 164 |
| 23 | 52,498,615 | 255,467 | 13 | 43 | 566 | 333 |
| 24 | 62,317,253 | 228,628 | 12 | 0 | 83 | 107 |
| 25 | 42,350,435 | 142,293 | 7 | 0 | 225 | 43 |
| 26 | 51,992,305 | 177,147 | 9 | 1 | 133 | 120 |
| 27 | 45,612,108 | 175,471 | 9 | 0 | 99 | 87 |
| 28 | 45,940,150 | 179,521 | 9 | 0 | 92 | 99 |
| 29 | 51,098,607 | 172,660 | 9 | 0 | 306 | 358 |

| Total | 2,489,385,779 | 8,234,367 | 418 | 574 | 7,764 | 6,298 |
|---|---|---|---|---|---|---|

757 Legend: BTA = chromosome number. GCR = genotype criteria sites. NoRS9 = number of site with
758 no RS number and with at least one alaternate allele in all 9 genotypes.
759

760 **Table 2     The overlap between the four methods for locating a likely novel variant site**

761

| Method | Genotype criteria (GCR) | Autozygosity by difference (ABD) | High-impact SIFT score (SIFT) | No registered SNP number (NoRS9) |
|---|---|---|---|---|
| GCR | **574 (12)** | | | |
| ABD | 12 (12) | **838 (575)** | | |
| SIFT | 1 (1) | 12 (12) | **7,764 (12)** | |
| NoRS9 | 4 (4) | 37 (36) | 8 (1) | **6,298 (36)** |
| GCR+ABD | | | 1 (1) | 1 (1) |
| GCR+SIFT | | | | 1 (1) |
| ABD+SIFT | | | | 1 (1) |
| GCR+ABD+SIFT | | | | 1 (1) |

762 Legend: The table shows the number of SNV in the final VCF file identified by each method.
763 (numbers in the BTA4 high-ABD region shown in parentheses).

764

765 **Table 3     Animal status by genotype for the 41 Sanger-sequenced animals at Chr4:**
766 **g.77173487A>T (ARS-UCD1.2 (GCF_002263795.1)).**

767

| Animal status | AA | AT | TT | Total |
|---|---|---|---|---|
| Calves | 4 | 2 | 7 | 13 |
| Known adult carriers (live) | 0 | 6 | 0 | 6 |
| Status unknown adults (live) | 13 | 9 | 0 | 22 |
| Total | 17 | 17 | 7 | 41 |

768

769 **Figure 1     A chromosome containing a novel autosomal recessive variant demonstrating**
770 **both the GCR and ABD methods**

771 Legend. Heat map of 28 controls (left side of the figure) and 7 cases (right side of the figure) for a
772 hypothetical chromosome containing a novel autosomal recessive variant. Chromosomes for the
773 individuals run from top to bottom of the figure. The colors represent homozygous major allele
774 genotypes (red), homozygous minor allele genotypes (green) and heterozygote genotypes (yellow).
775 Biallelic variants are assumed. The solid black bar across the controls represents the position of the
776 new variant with a GCR of all heterozygotes in controls and all the same homozygote in cases. The

777    ROH around the candidate position for cases are shown as a black rectangle running up and down the
778    chromosome until a heterozygote is found. The ABD method takes the mean length, in Kb, of ROH
779    in cases minus that in controls at each position to calculate the ABD score.

780    **Figure 2        Manhattan plots of the ABD analysis of nine WGS animals (Kb)**

781    Legend: P < 0.01 at ABD score = 9,034 Kb).

782

783    **Figure 3        Manhattan plot of the ABD analysis of the SNP-chip analysis based on the**
784    **genotypes found in the PCR analysis (Kb)**

785    Legend: These results were based on animals with phenotyping informed by the PCR results. (P <
786    0.001 at ABD score = 7,023 Kb).

787