



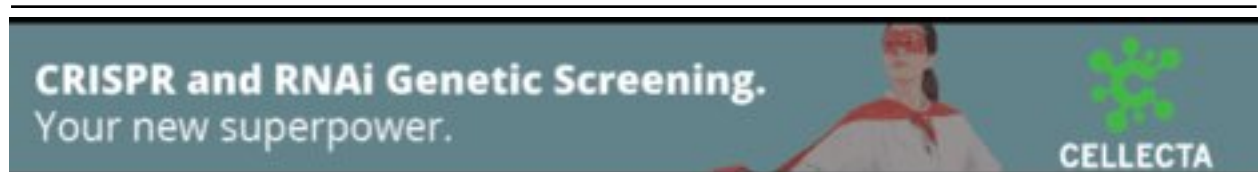
SNP-based quantitative deconvolution of biological mixtures: application to the detection of cows with subclinical mastitis by whole genome sequencing of tank milk

Wouter Coppieters, Latifa Karim and Michel Georges

Genome Res. published online June 26, 2020

Access the most recent version at doi:[10.1101/gr.256172.119](https://doi.org/10.1101/gr.256172.119)

P<P	Published online June 26, 2020 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see http://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **SNP-based quantitative deconvolution of biological mixtures: application to the detection of cows**
2 **with subclinical mastitis by whole genome sequencing of tank milk.**

3

4 *Wouter Coppieters¹, Latifa Karim¹, Michel Georges².*

5

6 ¹Genomics Platform, GIGA Institute, University of Liège. ²Unit of Animal Genomics, GIGA Institute &
7 Faculty of Veterinary Medicine, University of Liège.

8

9 Correspondence: michel.georges@uliege.be

10

11 **Abstract**

12 Biological products of importance in food (f.i. milk) and medical (f.i. donor blood derived products)
13 sciences often correspond to mixtures of samples contributed by multiple individuals. Identifying
14 which individuals contributed to the mixture and in what proportions may be of interest in several
15 circumstances. We herein present a method that allows to do this by shallow whole genome
16 sequencing of the DNA in mixed samples from hundreds of donors. We demonstrate the efficacy of
17 the approach for the detection of cows with subclinical mastitis by analysis of farms' tank mixtures
18 containing milk from as many as 500 cows.

19

20 **Introduction**

21 Mastitis, i.e. the inflammation of the udder, is the most important health issue in dairy cattle. It is
22 estimated to cost European farmers > 1 billion € per year in treatment and milk loss (Hogeveen et al.
23 2011). Upon inflammation, immune cells migrate in the udder and milk. While milk from healthy
24 cows typically contains > 100,000 cells per milliliter (ml) of milk, these numbers (referred to as Somatic
25 Cell Counts or SCC) typically increase into the millions in case of mastitis. Prior to the manifestation
26 of overt clinical symptoms, SCC progressively increase in the milk of cows developing mastitis: SCC ≥
27 200,000 are typically considered to be a sign of pre- or sub-clinical mastitis. Both yield and quality of
28 the milk of cows with subclinical mastitis is reduced (Schukken et al., 2003). Mastitis is routinely
29 managed by periodically counting SCC in milk samples to preemptively identify cows developing
30 subclinical udder inflammation. As profit margins decrease, farmers tend to forgo milk testing thereby
31 compromising health management. Cost-effective alternatives for rapid detection of cows with
32 subclinical mastitis are hence needed (Viguier et al., 2009).

33 The milk obtained from individual cows is typically collected in one or more large "milk tanks" on the
34 farm, before being shipped to dairy factories. We previously proposed that somatic cell counts (SCC)

35 in the milk of individual cows could be estimated by measuring the allelic frequencies in the tank milk
36 for sufficient numbers of SNPs, provided that all cows contributing milk to the tank be genotyped once
37 for the corresponding variants. Thus, the proposed method would allow the identification of a
38 minority of cows with subclinical mastitis by regularly analyzing a single sample containing a mixture
39 of milk from all the cows on the farm, hence considerably reducing costs (Blard et al. 2012). Prior to
40 ~2010 estimation of breeding values to select the best dairy sires and dams used pedigree-based
41 estimates of kinship. Since then, selection methods increasingly use genome-wide SNP information
42 in a process referred to as “genomic selection” (GS) (Georges et al. 2019). As GS is becoming routine
43 in dairy cattle (including for dams), herds that are fully genotyped with genome-wide SNP arrays are
44 becoming standard, and the proposed method feasible. However, as GS typically relies on the use of
45 low-density SNP arrays, the basic method proposed in Blard et al. (2012) is only effective for small
46 farms (≤ 100 cows). We herein demonstrate that by combining low density SNP genotyping or
47 shallow sequencing of the cows and tank milk’s DNA with in silico genotype imputation, individual SCC
48 can be accurately determined and cows with subclinical mastitis effectively identified even in the
49 largest farms (≥ 500). The proposed method has the potential to improve the monitoring of udder
50 health in dairy farms, and to allow the tracing of the origin of bulk animal food products other than
51 milk.

52

53 Results

54 **Principle of the proposed method.** Assume that cows and tank (i.e. the reservoir in which the milk of
55 the cows is collected) milk are genotyped for a collection of SNPs. Assume that the interrogated SNPs
56 are biallelic, each characterized by a *A* (say the allele of the reference genome) and a *B* allele (say the
57 alternate allele). If all cows contribute identical amounts of DNA to the milk, the expected proportion
58 of the *B* allele (commonly referred to as “B-allele frequency” when analyzing SNP array data
59 particularly to search for Copy Number Variants) in the tank milk corresponds to the frequency of the
60 *B* allele in the farm’s cow population. The actual DNA amount contributed by each cow depends on
61 the volume of milk that she produced and its SCC. Unequal DNA contributions will cause slight
62 departures from the expected *B* allele frequencies in the tank milk. Integrating these shifts over a
63 large number of SNPs in conjunction with the known genotypes of individual cows allows for the
64 estimation of the relative DNA contribution of each cow. This can for instance be achieved using a set
65 of *m* linear equations in which the “B-allele frequency” of each SNP *j* (of *m*) is modelled as the sum
66 (over *n* cows) of the products of the dosage of the *B* allele in the genotype of cow *j* (d_{ij} , known from
67 her SNP genotype) multiplied by the proportion of DNA contributed by cow *i* (f_i) to the milk. The
68 proportions of DNA contributed by each cow can then be estimated using for instance least square

69 methods. Accounting for individual milk volumes and for the SCC in the tank milk allows for the
70 estimation of SCC for individual cows (Fig. 1 and Methods).

71 **Evaluating the proposed method by simulation.** We first evaluated the proposed method by
72 simulation (cfr. Methods). Genotyping the cows and the tank milk using 10K SNP arrays (i.e. low-
73 density (LD) arrays as generally used in the context of genomic selection) allowed for the accurate
74 estimation of individual SCC for farms with up to 100 cows ($r \geq 0.9$, where r is the correlation
75 between real and estimated SCC) (scheme I). However, farms with > 100 cows are increasingly
76 common. Medium- (MD, f.i. 50K) and high-density (HD, f.i. 700K) SNP arrays would be needed for
77 the approach to be effective in farms with ≥ 250 or ≥ 500 cows, respectively. Yet – being too
78 expensive - this is presently not a viable proposition (Fig. 2A and Supplemental Table 1). We therefore
79 envisaged a second scheme (II) in which the cows would still be genotyped with LD SNP arrays (as
80 done in practice) yet imputed (Marchini & Howie 2010) to whole genome (8 million SNPs in the
81 simulations) using a sequenced reference population (f.i. Daetwyler et al. 2014), while the DNA of the
82 tank milk would be genotyped by shallow whole-genome sequencing (SWGS). We found that under
83 this scenario sequencing the tank milk at a depth of $0.25\times$ was sufficient for farms with 100 cows,
84 $0.5\times$ for farms with 250 cows, and $2\times$ for farms with 500 cows (Fig. 2B). Accuracies were not
85 significantly affected by the density of the SNP arrays, i.e. the method performed as well with LD as
86 with MD arrays (Suppl. Fig. 1). Anticipating further advances in sequencing technology, we also
87 envisaged a scheme (III) in which both cows and tank milk would be genotyped by SWGS. We found
88 that a $1\times$ sequencing depth of the tank milk would be sufficient when combined with a $0.25\times$ depth
89 for 100 cows, while a $5\times$ sequencing depth of the tank milk would be needed in combination with
90 $0.25\times$ depth for 250 cows and $1\times$ depth for 500 cows (Fig. 2C&D). In scheme III, allelic dosage in the
91 cows is directly measured from the number of alternative and reference alleles in the sequence reads.
92 We further explored the effectiveness of augmenting the cow genotype information from SWGS by
93 imputation (scheme IV). This proved to be effective, reducing the required sequence depth to $0.25\times$
94 for tank milk and $0.25\times$ for 100 cows, to $1\times$ for tank milk and $0.25\times$ for 250 cows, and to $5\times$ for tank
95 milk and $0.25\times$ for 500 cows (Fig. 2). The previous simulations make a number of assumptions that
96 may not apply in the real world: (i) SNPs were sampled from a uniform distribution (i.e. rare and
97 common SNPs equally represented), (ii) SNPs were assumed to be in linkage equilibrium, (iii) cows on
98 the farm were assumed to be unrelated, and (iv) milk volumes were assumed to be known without
99 error. To more accurately mimic real conditions we repeated the simulations by (i) sampling
100 genotypes from a phased dataset of 750 Holstein-Friesian whole genome sequences (hence properly
101 accounting for true MAF distribution, true linkage disequilibrium (LD) and relatedness - many of the
102 sequenced animals were related as parent offspring, full- or half-sibs), and (ii) adding a normally

103 distributed error with mean 0 and standard deviation of five liter to the simulated milk volumes
104 (assumed to be normally distributed with mean of 30 liter and standard deviation of 10 liter). This
105 error rate corresponds approximately to that expected when having to estimate the daily milk volume
106 from the total lactation yield using a standard lactation curve (Atashi 2019). We assumed in these
107 simulations that the genotypes of the cows were known without error and that the milk was
108 sequenced at a depth ranging from 0.25× to 5× as before. MAF, LD and relatedness jointly had a
109 relatively modest impact on the accuracy of the method, which could be compensated for by
110 increasing the sequencing depth of the milk to five-fold and still allowing for accurate estimates even
111 in farms with 500 cows. Estimating the milk volume with error had a more pronounced impact on
112 the accuracy making it difficult to reach a correlation reaching 0.9 in farms with 500 cows (Fig. 2).

113

114 ***Real-world application of the proposed method.*** To test the feasibility of our method in the real
115 world, we first collected cow (blood) and tank (milk) samples from a farm milking 133 Holstein-Friesian
116 cows. When only using genotypes from the Illumina LD arrays (17K SNPs) for both cows and tank milk
117 (scheme I), correlations between predicted and measured SCC were 0.91 (or 0.79 when ignoring one
118 cow with SCC > 3 million). We then imputed the cows to whole genome (13M SNPs) using a reference
119 population of ~750 whole genome sequenced Holstein-Friesian animals, and sequenced the tank milk
120 at ~3.5× depth. The corresponding correlations (scheme II) were 0.97 (0.95) when using all sequence
121 information, or 0.96 (0.92) when down-sampling sequence information as low as 0.1× depth (Fig. 3A).
122 We next performed a similar experiment on a farm milking 520 Holstein-Friesian cows. The
123 correlation between predicted and measured SCC was 0.78 (or 0.42 when ignoring 23 cows with SCC
124 > 3 million) when only using information from the LD array for both cows and tank milk (scheme I).
125 When imputing the cows to whole genome (13M SNPs) and sequencing the milk at ~3.5× depth
126 (scheme II), the correlation increased to 0.89 (0.83). Down-sampling the sequence information to
127 0.1× depth reduced the correlation to 0.79 (0.57) (Fig. 3B).

128 As shown in both farms, correlation estimates are affected by SCC spread: small numbers of cows with
129 very high SCC tend to inflate r . We therefore computed accuracies, computed as the proportion of
130 correctly classified cows for different SCC thresholds, which is how farmers would likely use the
131 information. It can be seen that for a threshold value of for example 500,000 SCC, accuracies > 0.85
132 were obtained when sequencing (scheme II) the tank milk at respectively 0.1× (133 cows) and 3.5×
133 depth (520 cows). Thus - as predicted by the simulations - scheme I provided adequate precision for
134 the farm with 133 cows, but not for the farm with 520 cows. However, in this large farm, combining
135 SWGS of the tank milk with whole genome imputation of the cows (i.e. scheme II) was indeed effective
136 (Fig. 3).

137 As costs per bp continue to decline, sequencing is likely to replace array-based genotyping in the
138 future. To test the feasibility of schemes III and IV (i.e. genotype the cows by SWGS rather than with
139 SNP arrays, without (III) and with (IV) imputation), we collected samples from a farm with 120
140 Holstein-Friesian cows. All cows were genotyped with the Illumina LD array (17K) as well as sequenced
141 at average 1.08× depth (range: 0.26-1.73). The milk was sequenced at ~3.5× depth. The correlation
142 between predicted and measured SCC was 0.97 (or 0.96 when ignoring one cow with SCC > 3 million)
143 under scheme I. Under scheme III, correlations were 0.82 (0.83) when sequencing the milk at 3.5×
144 and 0.75 (0.76) when down-sampling the milk to 0.1×. We then imputed the sequenced cows to HD
145 (770K SNPs) using a population of 800 reference animals genotyped with the HD array (scheme IV).
146 The correlation increased to 0.93 (0.94) when sequencing the milk at 3.5× and to 0.83 (0.77) when
147 down-sampling the milk to 0.1× (Fig. 3C). Accuracies at SCC threshold of 500,000 were 0.96 (scheme
148 I), 0.95 (3.5×) and 0.80 (0.1×) (scheme II), 0.82 (3.5×) and 0.81 (0.1×) (scheme III), and 0.95 (3.5×)
149 and 0.88 (0.1×) (scheme IV) (Fig. 3C). In summary, (i) combining cow genotyping using SNP arrays
150 with genome-wide imputation with SWGS of tank milk allows for cost-effective identification of cows
151 with subclinical mastitis even in farms with as many as 500 cows per milk tank, and (ii) as sequencing
152 costs continue to decline, arrays-based targeted SNP genotyping of the cows could be replaced by
153 genotyping by SWGS and yield comparable results.

154 **Monitoring SCC dynamics with the proposed method.** Farmers typically measure individual SCC once
155 a month or less. Yet, SCC may rapidly change. The SCC measured on the milk testing date may not be
156 a reliable indicator of the cow's udder health during the intervening period. To examine the SCC
157 dynamics over time, we collected 20 tank milk samples over a 100-day period (day -84 to +17 from
158 day of milk testing) for the farm with 120 cows. Milk samples were genotyped using the Illumina LD
159 array, and individual SCC estimated using scheme I. Fig. 4A shows the SCC predicted every 5 days on
160 average for the 120 cows, sorted by SCC measured on day 0 (=milk testing day). Of note, the
161 correlation between the SCC measured on day 0 and the average of the SCC estimates for the 21
162 collection dates was low ($r = 0.52$)(Fig. 4B) and decreased rapidly with the number of days from milk
163 testing day (Fig. 4C).

164

165 Discussion

166 We herein demonstrate that by combining array-based SNP genotyping and whole-genome
167 imputation for the cows with SWGS of the tank milk, it is possible to accurately estimate SCC for
168 individual cows and hence effectively identify animals with subclinical mastitis even for tanks
169 collecting milk for >500 cows, and this by performing a single analysis for the entire herd. Reagent
170 costs to sequence a mammalian genome at 1-fold depth are now <20€ thus making this a cost-

171 effective proposition. As a matter of fact, the method is being deployed in the field in several
172 countries.

173 Implementing the method requires all cows on the farm to be genotyped. This will increasingly
174 correspond to reality as genotyping costs continue to decrease and genomic selection is more and
175 more used for the selection of cows. In 2016 more than 1.2 million dairy cows had been reportedly
176 genotyped in the US alone (Wiggans et al. 2017) and present worldwide numbers are likely ≥ 3 million.
177 In addition, a reference population of a few hundred animals of the breed of interest that are either
178 HD genotyped (700K) or better whole-genome sequenced are required for accurate imputation. Such
179 reference populations are already available for the most important dairy cattle breeds (Daetwyler et
180 al. 2014; Charlier et al. 2016), and could be easily generated for the remaining ones.

181 We show that SCC are dynamic and rapidly change over time. SCC measured on day 0 are poor
182 indicators of SCC in previous and future weeks: cows with high SCC on the day of milk testing may
183 have low SCC a few days later (or earlier) and vice versa. The proposed method would allow tighter
184 monitoring of SCC hence improving udder health management. More frequent monitoring of SCC for
185 large number of cows may reveal interindividual differences with regards to SCC dynamics that may
186 be correlated with mastitis resistance, heritable and hence amenable to selection including by GS.

187 Sequencing of the DNA in the tank milk allows simultaneous characterization of the tank's
188 microbiome. As a matter of fact, $\sim 1\%$ of reads in this study mapped to bacterial genomes. This
189 information may be very useful both from a farm health management point of view as well as from a
190 downstream dairy processing point of view. Whole genome sequence data of bulk milk also informs
191 about the herd frequency of functional variants such casein variants affecting consumer health or
192 processing properties (Brooke-Taylor et al. 2017), or variants causing inherited defects or embryonic
193 lethality in cows (Georges et al. 2019). In many countries, it is not allowed to add milk from cows
194 being treated with antibiotics to the tank. As suggested before, the proposed approach can be
195 adapted to verify whether a specific cow did contribute milk to the tank or not (f.i. by testing the
196 significance of the corresponding cow effect in the linear model) (Blard et al. 2012). The described
197 method may have applications in tracing the origins of bulk animal food products other than milk, as
198 well as in monitoring the composition of mixed-donor blood-derived transfusion products.

199

200 **Methods**

201 **Simulated data.** Reference scheme (I): We simulated farms with n (25, 50, 100, 250 and 500) cows
202 contributing milk to the tank. Cows were genotyped with SNP arrays for m (10K, 50K, or 750K) markers
203 without error. Minor Allele Frequencies (MAFs) were sampled from a uniform $]0,0.5]$ distribution, and
204 genotypes from the corresponding Hardy-Weinberg distributions. SCC of individual cows (SCC_i) were

205 simulated by sampling values from a Weibull distribution with scale parameter $\alpha=1$ and shape
 206 parameter $\beta=2$, and multiplying the ensuing value by 200,000. Exact B-allele frequencies of individual
 207 SNPs (BAF_j) in the milk were determined for each SNP j based on the combination of cellular
 208 contribution of the n cows to the milk, and their genotype. It was assumed that B-allele frequencies
 209 were estimated with a normally distributed error $N(0, 0.0025)$ (i.e. SE = 0.05), yielding $m \widehat{BAF}_j$.
 210 Scheme II: Same setting as in the reference scheme with the following additions. For cows genotyped
 211 for 10K or 50K SNPs, we simulated imputation by augmenting the data to 8 million (M) genotypes
 212 using an error model mimicking real, MAF-dependent imputation accuracy. The error model was
 213 constructed using a real data set for 800 unrelated Holstein-Friesian individuals that were genotyped
 214 for the Illumina 777K array. This data set was split into a set of 200 and a set of 600 individuals. The
 215 set of 200 was reduced first to the genotypes interrogated by the Illumina 10K (LD) array and then to
 216 the genotypes interrogated by the Illumina 50K SNP arrays. The reduced SNP sets were imputed back
 217 to the content of the Illumina 777K (HD) SNP array using the 600 individuals as reference population.
 218 The frequencies of imputing a given genotype depending on the real genotype, were scored for MAF
 219 bins of 0.01 separately for the LD and 50K array data. We simulated genotyping-by-sequencing of tank
 220 milk as follows. For each of the 8M SNP positions, we sampled local read depth ($r \in \text{integers}$) from a
 221 Poisson distribution with mean C , where C is the average genome-wide coverage (0.25, 0.5, 1, 2 or 5).
 222 We then sampled r reads, each with a probability = BAF_j (computed as above) of being the B-allele.
 223 Scheme III: Individual SNP genotypes and tank B-allele frequencies (BAF_j) were generated as in
 224 scheme I (genotypes at 8 M SNP positions). It was assumed that milk tank was genotyped by SWGS
 225 at average coverage of C (0.25, 0.5, 1, 2 or 5) and cows were genotyped by SWGS at average coverage
 226 of C (0.25, 0.5, or 1). Genotyping-by-sequencing of individual cows was simulated by (i) sampling, for
 227 each of 8M SNP positions, local read depth ($r \in \text{integers}$) from a Poisson distribution with mean C ,
 228 and (ii) sampling r reads with probability 0, 0.5 or 1 to be the alternate allele (B) depending on the
 229 genotype of the cow (AA, AB or BB). Genotyping-by-sequencing of the tank milk was done as in
 230 Scheme I. Scheme IV: Identical to scheme III except that cow genotypes were generated at 8M SNP
 231 position using a MAF- and sequence-depth dependent imputation error model. The error model was
 232 constructed using available SWGS data down sampled to $1\times$ (176 cows) or $0.25\times$ coverage (192 cows).
 233 The cows were imputed to HD (777K SNPs) using a reference population of 800 unrelated Holstein-
 234 Friesian individuals that were genotyped with the Illumina 777K array. At each of the 777K SNP
 235 positions, the likelihood of the sequence data under the three possible genotypes (AA, AB and BB),
 236 were computed following Chan et al. (2016), as:

$$237 \quad L(nr_A, nr_B | "AA", \varepsilon) = \binom{nr_A + nr_B}{nr_B} \times (1 - \varepsilon)^{nr_A} \times \varepsilon^{nr_B}$$

$$238 \quad L(nr_A, nr_B | "AB", \varepsilon) = \binom{nr_A + nr_B}{nr_B} \times 0.5^{(nr_A + nr_B)}$$

$$239 \quad L(nr_A, nr_B | "BB", \varepsilon) = \binom{nr_A + nr_B}{nr_B} \times (1 - \varepsilon)^{nr_B} \times \varepsilon^{nr_A}$$

240 where nr_A (respectively nr_B) is the number of A (respectively B reads) and ε is the sequencing error
 241 rate set at 0.01. The corresponding $\log_{10} L$ were used as input for Beagle4 (Browning & Browning
 242 2009). Variant positions without sequence coverage in any of the 176 (192) cows (hence not imputed
 243 by Beagle4) were dealt with in a second round of imputation using Beagle5 (Browning et al. 2018).
 244 The imputation accuracy was evaluated in 0.01 MAF-bins by comparing imputed and real genotypes
 245 at the ~17K variant positions interrogated by the Illumina LD array.

246 **Real data.** Data set 1: We obtained a sample of tank milk from a farm in France milking 133 Holstein-
 247 Friesian cows. All had been genotyped with an Illumina LD array interrogating 17K SNPs using standard
 248 procedures. For all cows, genotypes were imputed to whole genome using a reference population of
 249 743 Holstein-Friesian animals sequenced at average depth of 15× (range: 4-48) and the Beagle
 250 software (v5.0)(Browning & Browning, 2009) yielding allelic dosages for a total of 13 million SNPs.
 251 Individual milk records, including volume and SCC (cells/ml) measured on the day of the sample
 252 collection, were obtained for all cows that had contributed milk to the tank. DNA was isolated from
 253 1.5 ml tank milk using the NucleoMag kit (Macherey-Nagel). The tank milk DNA was first genotyped
 254 using the Illumina LD array interrogating 17K SNPs. An Illumina compatible NGS library was then
 255 prepared with 50ng of genomic DNA using the KAPA HyperPlus kit (Roche). Sequencing was
 256 performed on a NextSeq 500 instrument (Illumina), yielding 63 million paired end reads of 2×75 bp,
 257 corresponding to a genome coverage of 3.5×. Reads were mapped to the bosTau8 genome build using
 258 BWA-MEM (Li 2013). Reference (R) and alternate (A) alleles were counted at 13M SNP positions of
 259 the HD array using the Bam-ReadCount tool (<https://github.com/genome/bam-readcount.git>) for
 260 reads with a minimum mapping quality of 30. Data set 2: We obtained samples of tank milk from a
 261 Belgian farm including milk from 520 Holstein-Friesian cows. Milk volume and SCC (cells/ml)
 262 measured on the same day, were obtained for all cows that had contributed milk to the tank. All cows
 263 were genotyped with the Illumina LD array interrogating 17K SNPs using standard procedures, and
 264 imputed to whole genome using whole genome sequence data (average depth: 15×; range: 4×-48×)
 265 from 743 Holstein-Friesian animals as reference and the Beagle software (v5.0)(Browning& Browning
 266 2009) yielding allelic dosages for a total of 13 million SNPs. DNA extraction from the tank milk samples
 267 and genotyping with the Illumina LD (17K) array were conducted as for dataset 1. For sequencing of
 268 the tank milk, an illumina compatible sequencing library was prepared using 12 ng of DNA and the
 269 Riptide High Throughput Rapid Library Prep Kit (iGenomx). The library was sequenced on an Illumina
 270 NextSeq 500 2×150 paired end flow cell at 4× coverage. Data set 3: We obtained samples of tank

271 milk from a Belgian farm including milk from 120 Holstein-Friesian cows. Milk volume and SCC
 272 (cells/ml) measured on the same day, were obtained for all cows that had contributed milk to the
 273 tank. All cows were genotyped with the Illumina LD array interrogating 17K SNPs using standard
 274 procedures, and imputed to whole genome using whole genome sequence data (average depth: 15×;
 275 range: 4-48) from 743 Holstein-Friesian animals as reference and the Beagle software (v5.0)(Browning
 276 & Browning 2009) yielding allelic dosages for a total of 13 million SNPs. We additionally prepared
 277 Illumina compatible NGS library for each cow, using 12 ng of genomic DNA and the Riptide High
 278 Throughput Rapid Library Prep Kit (iGenomx). Libraries were sequenced on an Illumina NovaSeq S4
 279 2150 paired end flow cell at average 1.08× depth (range: 0.26-1.73). Cow genotype-by-sequencing
 280 data were imputed to HD (777K) density using a reference population of 800 Holstein-Friesian animals
 281 genotyped with the bovine HD Illumina array (777K SNPs) and the Beagle software (v5.0) (Browning
 282 & Browning 2009) yielding allelic dosages for a total of 777K SNPs. DNA extraction from the tank
 283 milk samples, genotyping with the Illumina LD (17K) array, and sequencing (coverage 4×) were
 284 conducted as for datasets 1&2. Data set 4: In addition to obtaining a sample of tank milk on the day
 285 of the milk recording (i.e. yielding the SCC measured using with a cell counter) for the Belgian farm
 286 with 120 cows, we weekly collected an additional 11 tank milk samples before and 9 samples after,
 287 spanning a total period of ~3 months. The corresponding DNA samples were genotyped using the
 288 Illumina LD (17K) array.

289 **Statistical model.** We defined a set of m linear equations of the form:

$$290 \quad \widehat{BAF}_j = \sum_{i=1}^n f_i \times d_{ij} + \varepsilon_j$$

291 in which f_i is the proportion of the DNA in the tank milk contributed by cow i , d_{ij} is the “dosage” of
 292 the alternate allele A for cow i and marker j , and ε_j is the error term for marker j . When genotyping
 293 the tank milk with arrays, \widehat{BAF}_j corresponds to the B-allele frequency estimated by Genome Studio
 294 (Illumina). When genotyping the tank milk by SWGS, \widehat{BAF}_j corresponds to the proportion of A reads
 295 at the corresponding genome position. For cow genotypes obtained with arrays, d_{ij} corresponds to
 296 0, 0.5 or 1 for genotypes AA, AB and BB, respectively. For cow genotypes obtained by imputation, d_{ij}
 297 is the dosage of the B allele estimated by Beagle. For cow genotypes obtained by SWGS, $d_{ij} =$
 298 $0.5 \times P("AB" | nr_A, nr_B, q_j) + P("BB" | nr_A, nr_B, q_j)$ where nr_A (respectively nr_B) is the number of A
 299 (respectively B reads) for marker j and cow i , and q_j is the population frequency of the B allele of
 300 marker j .

301

$$302 \quad P("AB" | nr_A, nr_B, q_j) = \frac{2q_j(1-q_j) \times 0.5^{nr_A} \times 0.5^{nr_B} \times \frac{(nr_A + nr_B)!}{nr_A! \times nr_B!}}{(1-q_j)^2 \times 1^{nr_A} \times 0^{nr_B} + 2q_j(1-q_j) \times 0.5^{nr_A} \times 0.5^{nr_B} \times \frac{(nr_A + nr_B)!}{nr_A! \times nr_B!} + q_j^2 \times 0^{nr_A} \times 1^{nr_B}}$$

303

$$304 \quad P("BB" | nr_A, nr_B, q_j) = \frac{q_j^2 \times 0^{nr_A} \times 1^{nr_B}}{(1 - q_j)^2 \times 1^{nr_A} \times 0^{nr_B} + 2q_j(1 - q_j) \times 0.5^{nr_A} \times 0.5^{nr_B} \times \frac{(nr_A + nr_B)!}{nr_A! \times nr_B!} + q_j^2 \times 0^{nr_A} \times 1^{nr_B}}$$

305

306 For SNPs j without usable information for cow i (f.i. genotyping failure or no covering reads) d_{ij} was
307 set at \widehat{BAF}_j .

308 The f_i 's were estimated by least square analysis, i.e. by minimizing $\sum_{j=1}^m \varepsilon_j^2$. When the tank milk was
309 genotyped by SWGS, we also performed a weighted least square analysis, i.e. we estimated f_i 's by
310 minimizing $\sum_{j=1}^m w_j \varepsilon_j^2$, where w_j is the coverage ($nr_A + nr_B$).

311 The SCC_i 's were calculated from the f_i 's

$$312 \quad SCC_i = SCC_{tank} \times V_{tank} \times f_i / V_i$$

313 Where V_{tank} and V_i are the volumes of milk in the tank and contributed by cow i , respectively.

314 The accuracies of the predictions were measured by the (i) correlation (r) between real and estimated
315 SCC_i , and/or (ii) the ability to discriminate animals with SCC above versus below a certain threshold
316 value measured as $(T_P + T_N)/n$, where T_P stands for the number of true positives, T_N for the number
317 of true negatives, and n for the total number of cows.

318 To test the effect of sequence depth on accuracy we sampled reads overlapping SNP positions with
319 probability x , such that $E(C \times x) = D$, where D is the desired sequence depth.

320

321 Data access

322 All sequence (fastq files) and genotype (vcf files) data used in this study have been submitted to the
323 European Nucleotide Archive (ENA)(<https://www.ebi.ac.uk/ena>) (accession number: PRJEB38123
324 /ERP121506) and European Variation Archive (EVA) (<https://www.ebi.ac.uk/eva/>)(accession number:
325 PRJEB38336). Additional information to rerun the analyses are provided as Supplemental table 2.

326

327 Acknowledgements

328 This work was funded by the Unit of Animal Genomics and by the ERC DAMONA grant to Michel
329 Georges. We are grateful to Jean-Bernard Davière, Pierre Lenormand, Bonny Van Ranst, Kristien
330 Neyens and Miel Hostens for providing the samples and information needed to conduct the
331 experiments. WC and MG designed experiments, analyzed data and wrote the manuscript. LK
332 performed experiments.

333

334 Competing interest statement

335 The proposed method is the subject of awarded (WO/2013/079289) and filed patents
336 ([PCT/EP2019/057628](#)).

337

338 **References**

339 Atashi H, Salavati M, De Koster J, Ehrlich J, Crowe M, Opsomer G, GplusE consortium, Hostens M. 2019.
340 Genome-wide association for milk production and lactation curve parameters in Holstein dairy cows.
341 *J Anim Breed Genet* doi:10.1111/jbg.12442.

342 Blard G, Zhang Z, Coppieters W, Georges M. 2012. Identifying cows with subclinical mastitis by bulk
343 SNP genotyping of tank milk. *J Dairy Sci* **95**:4109-4113.

344 Brooke-Taylor S, Dwyer K, Woodford K, Kost N. 2017. Systematic review of the gastrointestinal effects
345 of A1 compared with A2 β -casein. *Adv Nutr* **8**:739-748.

346 Browning BL, Browning SR. 2009. A unified approach to genotype imputation and haplotype phase
347 inference for large data sets of trios and unrelated individuals. *Am J Hum Genet* **84**: 210-223.

348 Browning BL, Zhou Y, Browning SR. 2018. A one-penny imputed genome from next generation
349 reference panels. *Am J Hum Genet* **103**:338-348.

350 Chan AW, Hamblin MT, Jannink J-L. 2016. Evaluating imputation algorithms for low depth genotyping-
351 by-sequencing (GBS) data. *PLoS ONE* **11**:e0160733.

352 Charlier C, Li W, Harland C, Littlejohn M, Coppieters W, Creagh F, Davis S, Druet T, Faux P, Guillaume
353 F, Karim L, Keehan M, Kadri NK, Tamma N, Spelman R, Georges M. 2016. Reverse genetic screen for
354 embryonic lethal mutations comprising fertility in cattle. *Genome Res* **26**: 1-9.

355 Daetwyler HD, Capitan A, Pausch H, Stothard P, van Binsbergen R, Brøndum RF, Liao X, Djari A,
356 Rodriguez SC, Grohs C, et al. 2014. Whole-genome sequencing of 234 bulls facilitates mapping of
357 monogenic and complex traits in cattle. *Nat Genet* **46**: 858-865.

358 Georges M, Charlier C, Hayes B. 2019. Harnessing genomic information for livestock improvement.
359 *Nat Rev Genet* **20**:135-156.

360 Hogeveen H, Huijps K, Lam TJGM. 2011. Economic aspects of mastitis: New developments. *N. Z. Vet.*
361 *J.* **59**:16–23.

362 Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
363 [arXiv:1303.3997](https://arxiv.org/abs/1303.3997).

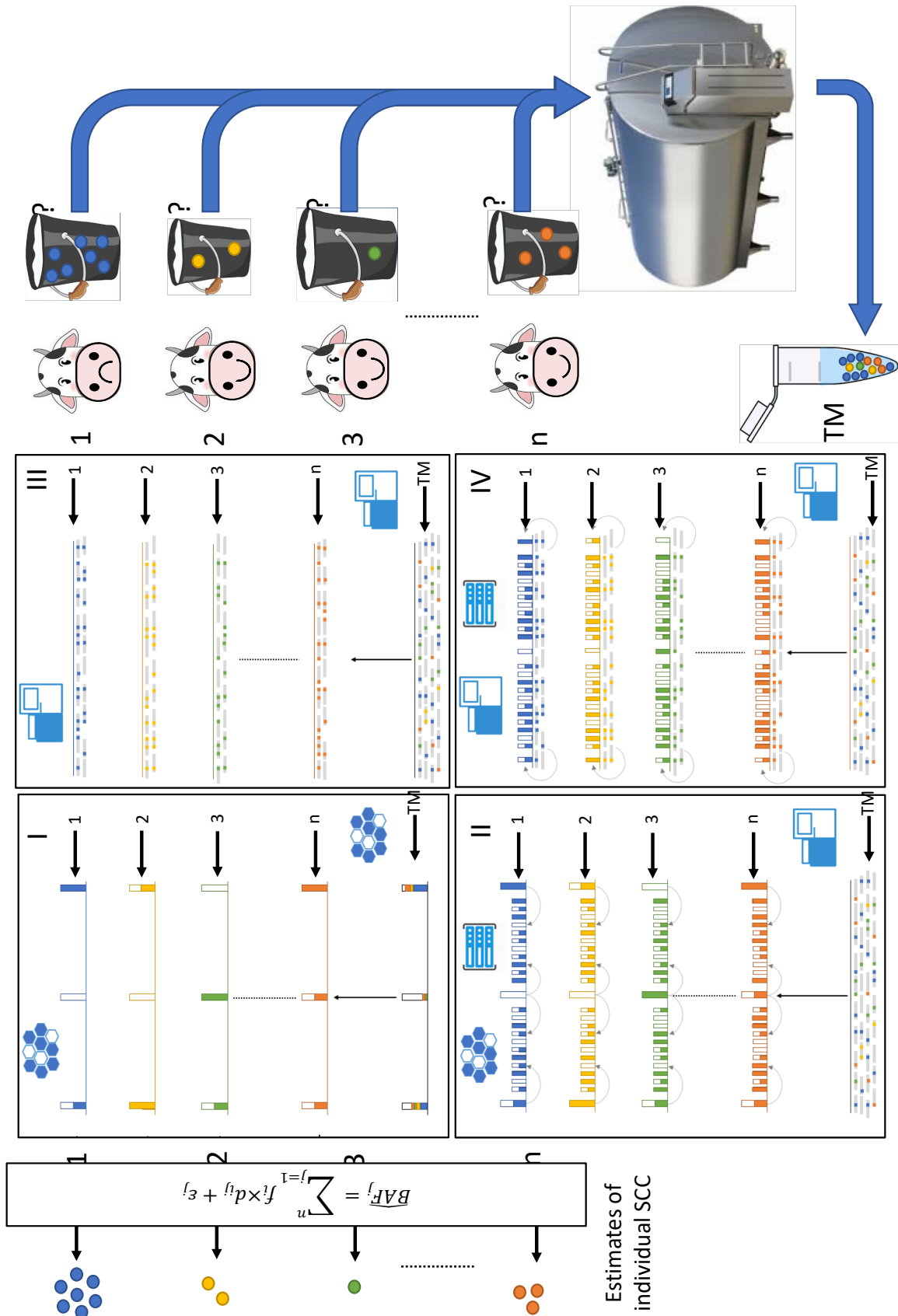
364 Marchini J, Howie B. 2010. Genotype imputation for genome-wide association studies. *Nat Rev Genet*
365 **11**:499-511.

366 Schukken YH, Wilson DJ, Welcome F, Garrison-Tikofsky L, Gonzales RN. 2003. Monitoring udder health
367 and milk quality using somatic cell counts. *Vet Res* **34**: 579-596.

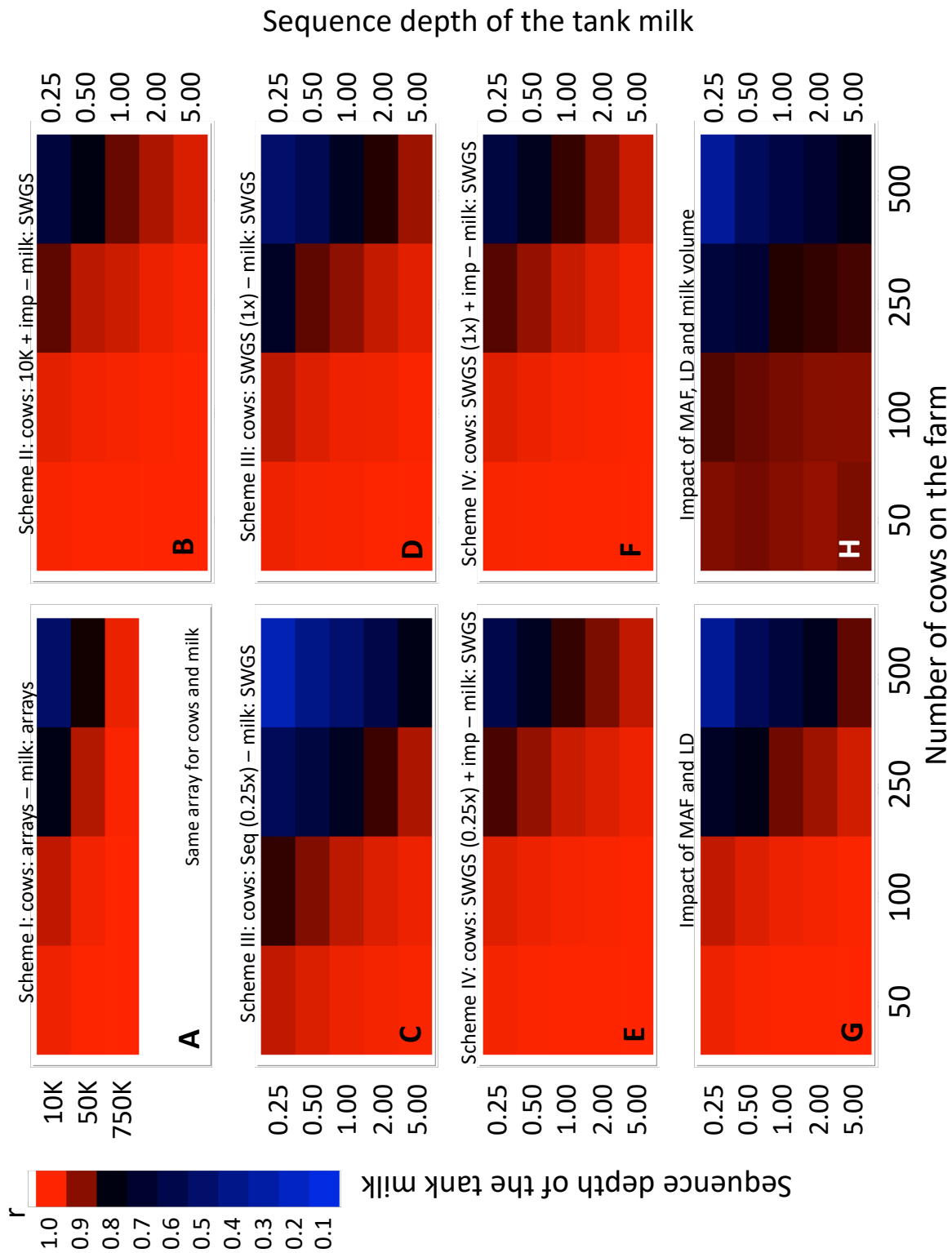
368 Viguier C, Arora S, Gilmartin N, Welbeck K, O’Kennedy R. 2009. Mastitis detection: current trends and
369 future perspectives. *Trends Biotechnol* **27**: 486-493.

370 Wiggans GR, Cole JB, Hubbard SM, Sonstegard TS. 2017. Genomic selection in dairy cattle: the USDA
371 experience. *Annu Rev Anim Biosci* **5**:309-327.

372 **Figure 1:** Estimating Somatic Cell Counts (SCC) in the milk of individual cows by analyzing a sample of
373 milk from the farm’s tank. Cows 1 to n contribute different amounts of milk (buckets of various sizes
374 in the figure) to the farm’s tank. The milk contains somatic cells (shown as small spheres in the milk
375 colored by cow) whose numbers reflect the health status of the cow’s udder. Cow 1 has higher SCC,
376 an indicator of subclinical mastitis. SCC are unknown upon milking (indicated by the “?”). Cows are
377 individually SNP genotyped once. In scheme I this is done using SNP arrays (illustrated by the mesh)
378 yielding genotype information for the limited number of interrogated SNPs (high bars) that can be
379 summarized by the B-allele frequency as shown (white: 0, halve colored: 0.5, fully colored: 1). SNP
380 genotypes of individual cows are coded in the same colors as the SCC. In scheme II, the genotypes of
381 the interrogated SNPs are augmented by imputation (illustrated by the computer rack), yielding
382 dosage information (B-allele frequency) for many more SNPs (small bars). In scheme III, cows are
383 genotyped individually by shallow whole genome sequencing (SWGS) (illustrated by the sequencer).
384 Sequence reads (gray lines) are aligned to the reference genome and alternate alleles at SNP positions
385 highlighted as color-coded tics. The B-allele frequency at specific SNP positions is measured as the
386 ratio of the number of reads with the B allele vs the total number of reads. In scheme IV, the genotype
387 information from SWGS is augmented by imputation improving the accuracy of the B-allele frequency
388 estimates for millions of SNPs (small bars). A small sample of milk (T(ank) M(ilk)) is periodically (f.i.
389 monthly or weekly) collected from the farm’s tank. DNA is extracted from TM and genotyped using
390 SNP arrays (scheme I) or SWGS (schemes II, III and IV). B-allele frequency for SNP j in the milk (\widehat{BAF}_j)
391 is estimated from the ratio of fluorescence intensities when using SNP arrays, or from the proportion
392 of reads with B allele in SWGS. The SCC of individual cows are estimated from a set of linear equations
393 modelling \widehat{BAF}_j as the sum of B allele dosage (d_{ij}) multiplied by the proportion of the DNA in the tank
394 contributed by cow i (f_i). The estimated proportions of DNA contributed by each cow correspond to
395 the values of f_i ’s that minimize the sum of squared errors (ϵ_j) over all SNPs. The SCC for individual
396 cows, *per se*, can be estimated as $SCC_i = SCC_{tank} \times V_{tank} \times f_i/V_i$, where SCC_{tank} is the SCC
397 measured in the farm’s tank, and V_i/V_{tank} is the proportion of the milk volume contributed by cow i .
398
399

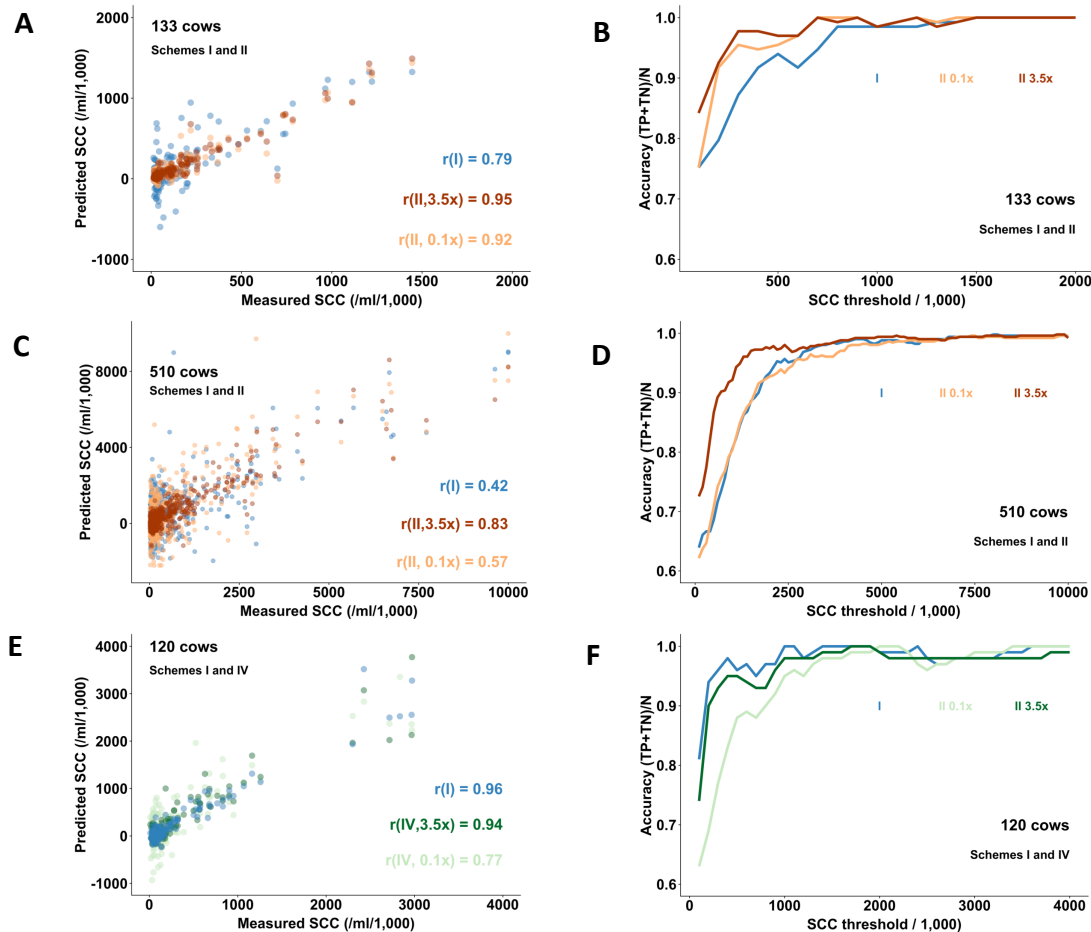


401 **Figure 2:** Evaluating the efficiency of the proposed approach by simulation. **(A)** Reference scheme I in
402 which individual cows and tank milk are genotyped with the same array interrogating 10K (LD), 50K
403 (MD) or 700F (HD) SNPs. **(B)** Scheme II in which individual cows are genotyped with a LD 10K SNP array
404 and imputed to whole-genome (8 million SNPs), while the tank milk is whole-genome sequenced at
405 depth ranging from 0.25× to 5×. **(C)** Scheme III in which individual cows (0.25×) and tank milk (range:
406 0.25× to 5×) are genotyped by shallow whole-genome sequencing (SWGS). **(D)** Same as C except
407 that individual cows are sequenced at 1× depth. **(E)** Scheme IV in which individual cows are genotyped
408 by SWGS (0.25×) followed by imputation to whole genome (8M SNPs), and tank milk is genotyped by
409 SWGS (range: 0.25× to 5×). **(F)** Same as E except that individual cows are sequenced at 1× depth.
410 **(G)** Scheme in which the cow genotypes are sampled from a real dataset hence conform to reality with
411 regards to distribution of MAF, LD and relatedness. Genotypes of the cows are assumed to be known
412 (very similar to II and IV) and tank milk genotyped by SWGS (range: 0.25× to 5×). **(H)** Same as G
413 except that the milk volume is estimated with error. The color code used to quantify the correlations
414 between predicted and real SCC is shown. Corresponding numerical values are provided in Suppl.
415 Table 1
416



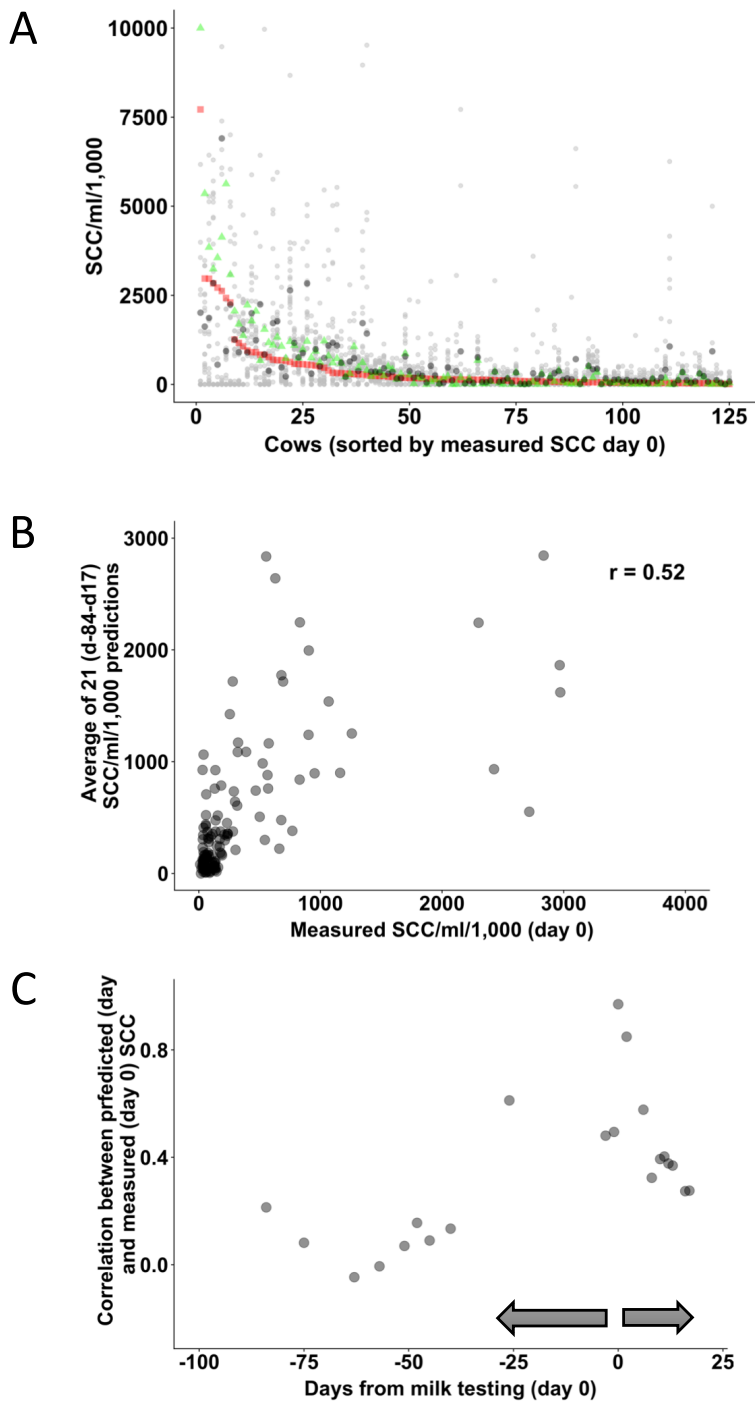
417
418
419
420

421 **Figure 3:** Correlation between predicted and measured SCC in the milk of individual cows (A,C,E), as
 422 as well as accuracies in classifying cows with SCC above and below a chosen threshold value (B,D,F), in
 423 farms with 133 (A,B), 520 (C,D) and 120 (E,F) cows, using scheme I (blue), scheme II (red), or scheme
 424 IV (green). Scheme I: cows and tank milk genotyped with LD SNP arrays (17K), no imputation. Scheme
 425 II: cows genotyped with LD array and imputed to 13M SNPs, tank milk sequenced 3.5x (red) or 0.1x
 426 (orange). Scheme IV: cows genotyped by whole-genome sequencing (1x) and imputation to HD, and
 427 tank milk sequenced at 3.5x (dark green) or 0.1x (light green).
 428



429
 430

431 **Figure 4: Evaluating SCC dynamics: (A)** SCC predicted using scheme A for 21 tank milk samples
432 collected over a 100-day period from 125 cows total. Small grey circles: 20 predictions per cow. Large
433 grey circles: average of 21 measurements per cow. Red square: SCC measured on day 0. Green
434 triangle: SCC predictions on day 0. **(B)** Relationship between SCC values measured on day 0 and
435 average of 21 predictions sampled over a 100-day period (days -84 to +17). **(C)** Correlations between
436 measured (day 0) and predicted (day x) SCC as a function of the number of days from day 0.



437