



D'Eath, R. (2012) Repeated locomotion scoring of a sow herd to measure lameness: consistency over time, the effect of sow characteristics and inter-observer reliability. *Animal Welfare*, 21:2, pp. 219-231. ISSN 0962-7286.

Copyright © 2012 Universities Federation for Animal Welfare.

This is the accepted version of the above article which has been published in final form at

<http://dx.doi.org/10.7120/09627286.21.2.219>

<http://hdl.handle.net/11262/8132>

Deposited on: 15 March 2016

1 **Repeated locomotion scoring of a sow herd to measure lameness: consistency over time, the effect of**
2 **sow characteristics and inter-observer reliability**

3 R.B. D'Eath

4 Animal Behaviour & Welfare, Sustainable Livestock Systems, SAC, West Mains Road, Edinburgh, EH9 3JG
5 rick.death@sac.ac.uk

6
7 Running head: Locomotion scoring of sows

8
9 **Abstract**

10 Investigating variability of scores between different observers, between animals and over time aids the
11 design of valid sampling methodologies for measuring animal welfare. Locomotion scores (0 to 5 scale)
12 were collected from i) 154 sows in one herd, using 3 to 5 (mean \pm s.d. = 3.95 ± 0.65) observers each
13 time, and scoring sows on up to 10 (mean \pm s.d. = 4.8 ± 2.7) occasions over a 19 month period, and ii) for
14 123 of these sows, locomotion scoring also took place prior to farrowing and at weaning.

15 The distribution of scores was highly skewed towards low scores (0: 84.8%, 1: 9.5%, 2: 4.0%, 3+: 1.7%).
16 Sows showed moderate consistency in their scores over time ($W=0.496$, $p<0.001$), and later parity sows
17 had higher scores ($\chi^2_1 = 3.98$, $p = 0.049$), but there was no effect of stage in the reproductive cycle (days
18 pregnant, pre-farrowing, post-weaning). This suggests that infrequent visits to a farm (e.g. annual) might
19 give a fair picture of the extent of lameness if a representative range of parities was sampled, although a
20 larger study of more farms would be required to investigate this.

21 The three different types of agreement between observers (absolute differences, matching and
22 association) were assessed as follows: i) Analysis of absolute differences between observers showed
23 that the farm manager scored lower than researchers/technicians ($S = 102.6$, $p < 0.001$). ii) Exact
24 matching approaches suggested fair ($\kappa = 0.443$) or good (PABAK = 0.692) agreement. Agreement was
25 poorest for mild gait abnormalities (score 1 "stiff"), and agreement improved if scores were combined
26 into 'sound' (0-1) and 'lame' (2-5) categories ($\kappa = 0.653$). iii) Measures of association suggested
27 moderate agreement (Kendall's $W = 0.692$, $p<0.001$). Inter-observer reliability improved over time until
28 the 5th scoring event. To improve inter-observer agreement, observer training/practice and the use of
29 fewer categories are recommended, and inter-observer agreement should be checked regularly.

30
31 **Keywords:** animal welfare, inter-observer agreement, lameness, locomotion scoring, pigs, sows

32

33 Introduction

34 On-farm measurement of animal health and welfare is an important and current issue, to meet
35 consumers' demands for demonstrably high standards of farm animal welfare (Blokhus et al 2003;
36 Blokhus et al 2008). Setting standards and inspecting to ensure that they are met is a goal of
37 government agencies (Gibbens 2008) and of voluntary farm assurance schemes such as the UK Red
38 Tractor scheme, free range, ecological and organic (Main et al 2003; Main et al 2007; Veissier et al
39 2008). Membership of such schemes generally depends on the producer meeting certain 'design criteria'
40 (Rushen & De Passillé 1992) relating to the housing and resources provided to animals (such as stocking
41 density, drinkers, substrates), management (e.g. weaning age in pigs, age at slaughter) and
42 administration (e.g. keeping accurate records of the use of drugs). Conformance with these criteria is
43 generally assessed in a visit which takes place about once a year and takes less than a day to complete
44 (Main et al 2007).

45 With a few exceptions, direct assessment of health and welfare by inspecting the animals themselves
46 (animal-based or 'performance criteria', Rushen & De Passillé 1992) has not formed part of these
47 schemes. Recently an EU funded project 'Welfare Quality®' (Blokhus et al 2003) developed a
48 comprehensive animal-based scoring system for on-farm assessment of animal health and welfare for
49 pigs (Welfare Quality® 2009) and other housed species. The measures adopted were assessed for
50 validity (does the indicator really measure what it should), repeatability (across observers), and
51 feasibility (can it be assessed quickly enough to be included in a short visit). This process has been
52 described in general terms but not in detail (Keeling et al 2009; Knierim & Winckler 2009).

53 Integration of multiple measures into an overall assessment is a difficult part of Welfare Quality
54 (Botreau et al 2007a; Botreau et al 2007b; Botreau et al 2007c; Botreau et al 2009; Knierim & Winckler
55 2009) and of similar schemes (e.g. Main et al 2007). Even before reaching this stage though, there are a
56 number of difficulties (Knierim & Winckler 2009). For any single measure, there are already practical
57 constraints: on-farm scoring of animal welfare involves a sampled subset of animals from each age class
58 and housing type, often by one trained observer in an annual visit of less than one day (Mullan et al
59 2009). Some researchers have attempted to assess the effect of such low 'sampling intensity' on
60 reliability of measures. These include studies of changes over time with repeated visits (Winckler et al
61 2007), the effect of sampling different numbers of animals at each visit (Mullan et al 2009; Main et al
62 2010), inter-observer reliability (Brenninkmeyer et al 2007; Bokkers et al 2009) and test-retest reliability
63 (O'Callaghan et al 2003; Flower & Weary 2006; Bokkers et al 2009). These questions are not merely
64 academic: animal-based scoring systems could, and perhaps should (e.g. FAWC 2008) become the basis
65 on which producers are deemed to pass or fail the criteria of an assurance scheme, so certain standards
66 of reliability must be reached (De Passille & Rushen 2005). Using locomotion scoring to detect lameness
67 in sows as an example of an animal-based scoring method, the present study investigated the extent to
68 which individual sows varied over time in their scores, and whether scores were affected by parity and
69 stage of pregnancy or veterinary treatment. Observers had different levels of experience with sows, and
70 were all initially naïve to the scoring system, so the extent to which experience with the system
71 improved consistency in the absence of training in these varied observers was also examined.

72 Lameness occurs in a variety of captive species, and is of welfare concern as the animal experiences
73 pain, discomfort and reduced mobility (Knowles et al 2008; Flower & Weary 2009). Lameness has
74 production costs too due to costs of treatment and is responsible for premature culling of around 7 – 11
75 % sows (Lucia et al 2000; Anil et al 2005; Engblom et al 2007; Anil et al 2009; Jensen et al 2010).

76 Although automated methods of assessing lameness have been investigated (Gonzalez et al 2008;
77 reviewed in dairy cattle by Flower & Weary 2009), simple human observer scoring systems are still the
78 main method used as they are relatively cheap, reliable and easy to apply on farm (Pigs, Main et al 2000;
79 KilBride et al 2009; KilBride et al 2010; Cattle, Winckler & Willen 2001; Flower & Weary 2006; Flower &
80 Weary 2009; Rutherford et al 2009; Chickens, Kestin et al 1992; Garner et al 2002). For sows in most
81 countries, confinement in stalls during pregnancy and crates during farrowing and lactation makes on-
82 farm assessment of lameness difficult, as altered posture is most likely not as sensitive a measure of
83 lameness as locomotion scoring (KilBride et al 2010). Confinement thus potentially obscures the extent
84 of the problem. The move to group housing following the EU ban on individual stalls during pregnancy
85 (Council Directive 2001/88/EC, 2001) in 2013 will both increase the need for good locomotion in sows
86 and make lame sows easier to identify.

87 In the present study, I applied locomotion scoring to group-housed sows (*Sus scrofa*) on one farm, using
88 multiple observers at each scoring event to measure inter-observer reliability. Agreement between
89 observers was assessed in terms of i) absolute differences, ii) exact matching and iii) association.
90 Although these are conceptually and statistically complementary approaches, they are rarely all used on
91 the same dataset (but see Kaler et al 2009). Because the same sows were scored on several occasions, in
92 addition I investigated consistency of scores over time, and the factors affecting locomotion scores such
93 as sow parity, stage in the reproductive cycle and the application of veterinary medicines

94 **Methods**

95 ***Animals and Housing***

96 154 Large White × Landrace sows at SAC's research pig unit were the subjects of this study. The study
97 animals included maiden gilts up to parity 7 sows (mean ± s.d. parity = 2.89 ± 1.72), and included dry
98 sows at various stages from soon after weaning (waiting to return to oestrous and be served) through to
99 heavily pregnant sows. They were housed in groups of 1-6 (mean ± s.d. = 4.5 ± 1.6) to a pen. Figure 1
100 shows the building where the sows were housed during the study. The pens (3.60m × 6.45m) were
101 concrete-floored, with an enclosed straw-bedded area at the rear (3.60m × 2.50m), walled with
102 concrete blocks, with a 2.0m wide opening onto a solid-floored central dunging passage (3.60m ×
103 1.95m), and an access passageway plus six individual feeding stalls side by side at the front (each 0.50m
104 wide, 1.80m long). Sows were fed on a rationed quantity of a commercial sow diet suitable for their
105 size/age and stage of gestation once a day (0800h). At each side of the pen was a barred gate across the
106 width of the dunging passage, which could be swung across to shut the sows into the bedded area
107 (Figure 1). Water was available in each pen via a nipple drinker mounted on one of the gates. The pens
108 were arranged in two back-to-back rows of 9 pens (108 sow places) and an automated natural
109 ventilation system set at 14.5 °C maintained a temperature of 8.1 °C min – 21.6 °C max, mean 16.2 ± 4.1
110 °C during the study. At one end of the building, there was an empty concrete-floored service pen, with a

111 small amount of straw and 3 cm deep sawdust (4.8m × 3.1m), used for artificial insemination, positioned
112 adjacent to where two 'teaser' boars were individually housed.

113 ***Locomotion Scoring***

114 Sows were locomotion scored in two contexts- 1) systematic dry sow herd scoring, where all sows
115 currently in the dry sow house were scored at one time and 2) scoring of sows due to farrow, or recently
116 weaned sows when they were being moved to or from their farrowing accommodation.

117 ***Dry Sow Herd Locomotion Scoring***

118 This took place on 11 occasions in total. The first two scoring events took place 6 months apart, but after
119 that they took place at intervals one to two months (min 27 – max 202 days apart, mean ± s.d. = 56.6 ±
120 53.9 days) over a period of 566 days between December 2008 and June 2010. There was turnover in the
121 breeding herd (old sows being culled, new gilts coming in) and a proportion of sows were in the
122 farrowing house at any one time. As such, between 47 and 76 (mean ± s.d. = 67.6 ± 8.4) sows were
123 scored at each scoring event, and individual sows were scored on 1- 10 occasions (mean ± s.d. = 4.79 ±
124 2.66).

125 Between 3 and 5 observers (mean ± s.d. = 3.95 ± 0.65) were present at each scoring event. Observers B
126 and C attended all 11 scoring events, A and D attended 9 events and observer E attended 4. Each
127 observer independently assigned each sow a locomotion score according to the system shown in table 1
128 and did not discuss or compare scores with other observers. The scoring system was simplified from the
129 system developed by Main et al. (2000) for growing pigs (see also KilBride et al 2009). Observers were
130 also free to record comments for every sow. For the dry sow herd scoring, observers consisted of a
131 scientist with practical pig experience over a period of 11 years, but now primarily desk-based (A; the
132 first author), whose interaction with these sows was limited to the scoring sessions only, another
133 scientist (B) with 8 ½ years of regular experience with sows in general, and running a project involving
134 these particular sows, the pig unit manager (C) with over 28 years experience working with all ages of
135 pigs every day, and two technicians (D) with 3 ½ years and (E) 18 years experience of regularly working
136 with pigs, including these sows. Before scoring began, the scoring system was explained and discussed
137 by the observers to ensure it was well understood, but no formal training was undertaken.

138 In the morning before scoring took place, sows were spray marked for ease of identification. Scoring
139 began at 1330h and took approximately 45 – 70 mins. First, all sows on one side of the dry sow house
140 (Figure 1) were shut into their bedded area, and then the group furthest from the service area was let
141 out of their pen and moved to the service pen. Sows either walked or ran, and some required
142 encouragement to walk from a stockperson walking behind them. During scoring, it sometimes became
143 apparent that sows were difficult to move. These were scored by all observers present and then sows
144 which were slow to move (scoring 4 Limping) were only moved a short way before being returned, while
145 those scoring 5 'Downer' were scored in their home pen. Once sows had been moved to the service pen,
146 they were shut in for 3 to 5 minutes, with one person who moved amongst the sows encouraging them
147 to move to facilitate scoring. Sows were then returned to their home pen, where they were no longer
148 shut into the bedded area. Then the next group was moved to the service pen in a similar way, and so
149 on for all the pens of sows until one side of the building was complete, then the procedure was repeated

150 for the other side of the building. While sows were moving to the service pen, in the service pen and
151 moving back, they were continually observed and scored by the observers. Scorers made sure to focus
152 on each individual sow for at least 10 strides of walking during this period and to give their score once
153 this was complete. Sows which were identified as scoring 3 or higher which had not previously been
154 identified during the course of normal husbandry were examined following scoring and subsequently
155 given appropriate veterinary treatment. This occurred on 13 occasions during the study. Depending on
156 the suspected cause of the lameness, sows were treated with different combinations of drugs. Eight
157 sows were given the non-steroidal anti-inflammatory drug Metacam (5ml) was usually given once, but
158 for two sows this was repeated daily for up to 7 days. A three day course of the antibiotic depocillin (8-
159 12.5ml depending on weight) was given to 8 sows and a broad spectrum antibiotic (Baytril; 8ml) was
160 also used on two occasions.

161 ***Farrowing Sow Locomotion Scoring***

162 For 123 of the sows, locomotion scoring took place prior to farrowing and after weaning. When they
163 reached 109 days after service, (mean \pm s.d. 4.38 ± 1.84 days before farrowing), they were moved out of
164 the dry sow house, across a concrete outside area, into a separate farrowing building and placed
165 individually in loose-housed (non-crate) farrowing pens (PigSAFE, Baxter et al 2011). They were
166 locomotion scored as they walked. The scoring was all done by observer B (a scientist). Once their
167 piglets were weaned (mean 26.8 ± 3.1 days after farrowing), sows were moved back to the dry sow
168 house, and were again locomotion scored as they walked between buildings.

169 ***Statistical Analysis***

170 The nature of the dry sow herd locomotion scoring data led to a number of challenges for analysis. The
171 data were incomplete because not every sow was present at each scoring event: sows were sometimes
172 in the farrowing house, and there was some turnover in the herd (culls and replacements) during the
173 study period. Also, not every observer was able to attend every scoring event. So there were a lot of
174 'missing' data where a particular sow/ observer/ scoring event combination did not occur.

175 Non-linear mixed models for ordinal data (using SAS; Gilmour et al 1987; Keen & Engel 1997) were fitted
176 for each event separately to assess inter-observer reliability. An underlying latent continuous variable
177 for locomotion was assumed of the form: $y_{ij} = \beta_0 + u_i + e_{ij}$ where u_i is a sow-level normal random effect
178 and e_{ij} are independent and identically distributed normal errors. y_{ij} was estimated from observer
179 scores. The 0-5 scoring scale corresponds on the latent variable scale to 6 intervals: $(-\infty, 0)$, $(0, I_1)$,
180 (I_1, I_1+I_2) , $(I_1+I_2, I_1+I_2+I_3)$, $(I_1+I_2+I_3, I_1+I_2+I_3+I_4)$ and $(I_1+I_2+I_3+I_4, \infty)$ where I_1 , I_1+I_2 , $I_1+I_2+I_3$,
181 $I_1+I_2+I_3+I_4$ are the thresholds for the categories. They were estimated when a threshold model for the
182 latent variable was fitted as a generalized linear mixed model (GLMM) using the NLMIXED procedure in
183 SAS 9.1 (SAS Institute Inc., SAS Campus Drive, Cary, North Carolina 27513). Inter-observer reliabilities
184 were derived from the variance components. In practice, scoring categories had to be merged on an
185 event-by-event basis to enable model fitting when there were very few animals on specific scores.

186 A different approach to dealing with the 'missing' data was used to look at consistency over time in sow
187 locomotion scores. A subset of the data was used: 81 sows which were scored on four or more
188 occasions by the four observers that did the most scoring (A-D), were identified. Kendall's coefficient of

189 concordance was calculated across the first four scoring events for each sow, using the mean of the four
190 observers' scores at each one. Note that the four specific scoring events used could vary between the
191 sows in this subset.

192 Another challenge for model fitting was that there were a high proportion of zero scores in the data. To
193 model the different factors affecting locomotion scores (sow parity and stage of pregnancy), the scores
194 were first re-coded into 0-1 data (0 → 0, 1-5 → 1) and a Generalized Linear Mixed Model was used to fit
195 a binomial model (with a logit transformation; Genstat 11th Edition), with sow as a random effect, and
196 observer, scoring event, sow parity and days until next farrowing (either actual or estimated from
197 service records) as fixed effects (with parity and days until next farrowing being fitted as covariates).
198 Note that non-pregnant sows and pregnant sows were treated similarly in this analysis: non pregnant
199 sows were simply those with greater than 113 days until next actual or predicted farrowing.

200 The effect of veterinary treatment on sow locomotion score was analysed by comparing scores for each
201 observer using paired Wilcoxon signed-ranks test. Scores from 0 to 90 (mean ± s.d. = 32.8 ± 24.6) days
202 before treatment began were compared with scores from 6 to 190 days after (mean ± s.d. = 57.3 ± 54.3).
203 One observer (B) was present for all 13 of the relevant scoring events, but other scorers had some
204 missing data (A = 8 events, C = 12, D = 10).

205 Various non-parametric methods were applied to measure inter-observer reliability. These require that
206 there was no missing data at all, so data from the 498 occasions when four of the observers (A, B, C & D)
207 were present was used (8 scoring events: 1, 3, 4, 6, 7, 8, 9 and 10; which included 137 sows). Using these
208 data, three aspects of agreement were considered:

209 1) Whether observers differed systematically in the absolute level of scores awarded. This was assessed
210 using Friedman for an overall comparison, and Sign tests and Wilcoxon signed rank tests (Minitab 15,
211 2006) to compare each pair of observers.

212 2) Whether scores from different observers matched exactly. Proportion of agreement and kappa were
213 calculated overall, by pairs of observers, and for different scores (Minitab 15, 2006). The prevalence
214 adjusted bias adjusted Kappa (PABAK, Byrt et al 1993) was also calculated as this is preferred by some
215 researchers (Brenninkmeyer et al 2007; Rutherford et al 2009). Weighted kappa (using linear
216 weightings) were also calculated (AgreeStat Excel workbook, Advanced Analytics, 2010). In interpreting
217 Kappa and PABAK, the scale suggested by Byrt (1996) was used: 0 or less no agreement, 0.01-0.20
218 poor..., 0.21-0.40 slight..., 0.41-0.60 fair..., 0.61-0.80 good..., 0.81-0.92 very good..., 0.93-1.00 excellent
219 agreement. Finally, kappa statistics were calculated following a re-coding of the 0-5 data to 0-1 data in
220 two different ways: Either all scores above zero were re-coded as 1 (0 → 0, 1-5 → 1), or 0 and 1 were re-
221 coded as 0, and higher scores as 1 (0 and 1 → 0, 2-5 → 1). This was done for comparison with other
222 researchers who have scored animals in two categories of 'sound' and 'lame' (e.g. Rutherford et al
223 2009). Two different approaches were chosen because of the difficulty in classifying score '1' animals
224 (see Results).

225 3) Whether scores from different observers were associated. This was assessed using Kendall's
226 coefficient of concordance (W; Genstat 11, VSN International Ltd, 2008) as a measure of overall

227 agreement, and Spearman's (ρ) and Kendall's (τ) rank correlation coefficients were calculated between
228 each pair of observers (Genstat 11, VSN International Ltd, 2008).

229 Finally, farrowing sow locomotion scores from before farrowing and after lactation were compared
230 using Wilcoxon-signed ranks test for paired data (after calculating differences, data were not normally
231 distributed).

232 ***Ethical Note***

233 This study was given ethical approval by SAC's Animal Experiments Committee. As detailed above, any
234 Sows which scored 3 or higher (3 Lameness, 4 Limping, 5 Downer) were given appropriate veterinary
235 treatment, while those scoring 1 or 2 were investigated and treated if necessary. Since sows were being
236 checked daily as part of their routine husbandry, it was rare for lameness problems to be newly
237 identified as part of the scoring process.

238 **Results**

239 ***Distribution of scores***

240 The frequency of scores given overall by each of the observers (i.e. for all sows over all scoring events) is
241 shown in figure 2. The vast majority of scoring events resulted in the observer giving the sow a score of
242 0: Normal (mean % of scores at each level: 0: 84.8%, 1: 9.5%, 2: 4.0%, 3: 1.2%, 4: 0.4%, 5: 0.2%).

243 ***Consistency of repeated scores from the same individual sows***

244 To estimate consistency over time of sow scores, the mean locomotion score of observers A-D, from 81
245 sows at four scoring events were used. These scores for this subset were distributed similarly to the
246 dataset as a whole as follows: 78.1% had mean score 0, 16.4% had a mean score between 0.25 and 1,
247 4.0% had a mean score between 1.25 and 2, 1.5% scored 2.25 or higher. There was a significant
248 moderate level of association between sows: Kendall's $W = 0.496$, $p < 0.001$, suggesting that individual
249 sows showed similar scores over repeated scoring events.

250 ***Factors affecting locomotion score***

251 GLMM models showed that there were differences in locomotion scores (re-coded to 0-1) due to
252 observer (Wald statistic $\chi^2_4 = 114.22$, $p < 0.001$) and scoring event: inspection of estimated means
253 showed that scores were lower at later scoring events ($\chi^2_{10} = 38.02$, $p < 0.001$). A histogram showing the
254 distribution of sow parities in the study is shown in figure 3. Older (higher parity) sows had higher
255 locomotion scores ($\chi^2_1 = 3.98$, $p = 0.049$). A histogram showing the distribution of stage of pregnancy is
256 shown in figure 4. There was no effect of stage of pregnancy ($\chi^2_1 = 0.10$, N.S.). Locomotion scores
257 observed before and after sows were given veterinary treatment for lameness did not differ for any
258 observer.

259 ***Agreement between observers- do they differ?***

260 From this section forward, all analyses (unless stated otherwise) use data from the 498 occasions when
261 observers A-D were all present, as the mainly non-parametric methods used require no missing data.
262 There were differences between the observers in the proportion of sows scored in each category.
263 Observer C (pig unit manager) recorded a higher proportion of zeroes than the other observers (Figure

264 2). A Friedman test comparing observers for the same sow on the same occasion showed a highly
265 significant difference between observers, and inspection of the sums of ranks showed this was because
266 observer C's scores were lower than those of the other three observers ($S = 102.6$, $p < 0.001$; sums of
267 ranks $A = 1252$, $B = 1299$, $C = 1143$, $D = 1286$). When pairs of values were compared with Sign and
268 Wilcoxon tests, these confirmed that observer C was scoring lower than the other three, and observer A
269 gave lower scores than B and D, who did not differ (Table 2).

270 ***Agreement between observers- do they match?***

271 Raw proportions of agreement both overall and between pairs of observers were high (Table 3). The
272 proportion of agreement was considerably higher for 0 than for other scores. This is best illustrated by
273 the proportion of occasions on which all 4 observers agreed to give the same score (first row of table 3).
274 The overall kappa statistic of 0.443 is at a level which suggests only 'fair' agreement (Byrt 1996),
275 although the PABAK statistic at 0.692 suggests 'good' agreement.

276 When broken down by scores (table 3), kappa and PABAK were noticeably lower for score 1, perhaps
277 reflecting the difficulty in identifying and agreeing on the threshold between 0: 'normal' and 1: 'stiff', a
278 score used for minor locomotor anomalies.

279 When broken down by pairs of observers (table 2), values of kappa (and of weighted kappa and kappa
280 following combining of categories) were considerably lower for pairings involving observer C (A-C, B-C
281 and C-D) than for the other three pairings. Pairings which include observer C suggest 'slight' agreement,
282 while pairings including observers A, B and D agreed at 'fair' or 'good' level. Contingency tables of the
283 frequency of scores given by observer pairs B and D (Table 4) and B and C (Table 5) are shown to
284 illustrate good and poor agreement respectively. For observers B and D, scores which differ by 2 are
285 rare, and the greatest number of discordant scores occurs at the 0 vs. 1 level. For table 5 (observers B
286 and C- poor agreement), there are more scores which differ by 2, 3 or even 4. Again, the greatest
287 discrepancy occurs at the 0 vs. 1 level.

288 Values for weighted kappa, which takes into account the size of the disagreement between observers
289 (table 2) were considerably higher than those of kappa, suggesting that disagreement by a small degree
290 (e.g. by one point of the scoring system) was more common than disagreement by a larger amount.

291 0-5 scores were converted into 0-1 scores (sound/lame) by combining categories in two different ways
292 (Table 2): i) Scores of 0 were classed as 'sound' and scores of 1-5 were classed as lame. Or ii) scores of 0
293 or 1 were classed as 'sound' and scores of 2-5 were classed as 'lame'. The Kappa was considerably lower
294 for the first of these methods than for the second, consistent with the suggestion that the observers
295 showed higher levels of agreement about the higher (2+) categories than they did about the distinction
296 between 'Normal (0)' and 'Stiff (1)' pigs.

297 ***Agreement between observers- are they associated?***

298 Using all available data, a non-linear mixed model for ordinal data was fitted for each scoring event, to
299 estimate inter-observer reliability. Across the 11 scoring events, reliability was moderate to high,
300 ranging from 0.552 to 0.879. It was evident that reliability improved over time (Figure 5), reaching a
301 plateau by about the 4th or 5th scoring event.

302 Using data from the 498 occasions when observers A-D were all present, Kendall's coefficient of
303 concordance (W) showed that there was moderate agreement between the observers. Levels of
304 spearman's ρ and kendall's τ (Table 2) for pairwise agreement were moderate to high. Values were
305 again lowest for pairs involving observer C, although the difference was less marked.

306 Because pairings involving observer C showed a lower level of agreement than other pairings, Measures
307 of agreement between the other observers, A, B and D were also calculated. Kendall's W was 0.784,
308 Fleiss's Kappa was 0.557, and weighted Kappa was 0.684. These values were all higher by around 0.1
309 than the overall agreement measures (Table 2).

310 ***Farrowing sow locomotion scoring***

311 The distribution of sows' locomotion scores before farrowing was 0: 83%, 1: 27%, 2: 12%, 3: 1%, and
312 after weaning was 0: 89%, 1: 25%, 2: 7%, 5: 2%. There were no differences between the locomotion
313 scores before farrowing and after weaning ($W=439$, n for test ignoring ties = 45, $p=0.379$).

314 **Discussion**

315 The right-skewed distribution of scores (figure 2) was as expected. This shape of distribution is common
316 in locomotion scoring studies (e.g. KilBride et al 2009). Levels of sow lameness in the present study were
317 comparable to previous reports: KilBride et al (2009) studied 88 UK herds and found that 14.4% of
318 pregnant gilts and 16.9% of pregnant sows had abnormal gait (score of 1 or higher), while 1.0 and 1.8 %
319 had minimal weight bearing on an affected limb (score 3 or higher). Respective figures for the present
320 study of 15.2% and 1.8%. Further studies are needed to assess the welfare significance of these scores
321 (Main et al 2000), but it seems likely that reduced weight bearing on the affected limb is indicative of
322 some discomfort/pain. Approaches to measure the welfare significance of locomotion scores include the
323 use of motivational measures (Weeks et al 2002) or analgesics (Rushen et al 2007; Flower et al 2008;
324 Danbury et al 2000), but these have not as yet been applied to pigs.

325 The structure of the dataset, with inevitable 'missing' data as sows moved through the system,
326 complicated the analysis of consistency of sow scores. Kendall's coefficient of concordance using the
327 average of four observers showed a moderate level of consistency. This analysis corresponds with
328 studies in dairy cows which reported a degree of consistency over time in locomotion scores
329 (O'Callaghan et al 2003; De Rosa et al 2003). It suggests that sow lameness is often a chronic problem.
330 Although veterinary treatments were only applied on 13 occasions during the study, there was no
331 evidence that scores were improved after treatment compared with before. This result should be
332 viewed with some caution: it was not the primary aim of the study, and sample size was very low. Also
333 there was a large and variable amount of time between the treatment and the before and after scores.
334 The treatments may well have worked well in the short or even long term, but new causes of lameness
335 may have occurred before the next scoring event.

336 Analysis of predictive variables found that locomotion score (analysed as binomial data, 0 vs 1-5) was
337 affected by parity but not stage of pregnancy. Parity may have affected locomotion score because
338 heavier animals are more likely to become lame (because of the greater pressure on their feet and
339 joints), or because older sows have had longer to pick up an accidental injury which may then take some
340 time to resolve. In dairy cow studies, size, conformation and udder fill affect how cows walk (reviewed

341 by Flower & Weary 2009), and weight affects broiler locomotion scores (Kestin et al 1992). During this
342 study, observers commented that they had had difficulty with the lowest end of the scoring system. It
343 was felt that the system failed to reflect the diversity of 'normal' and 'abnormal' gaits, particularly in
344 sows differing in age, weight or stages of pregnancy. For example 'swagger of rear end while walking' (2-
345 slight lameness) was quite pronounced in some otherwise normal older sows, so the effect of parity on
346 locomotion score may partly reflect this change in gait of older sows, as any deviation from score 0
347 would have affected this analysis.

348 In terms of recommendations for welfare assessment study design and sampling intensity, the variation
349 across parities suggests that the sampling strategy of larger on-farm studies should take this into
350 account, ensuring that a representative cross-section of the range of parities on the farm is sampled.
351 The moderate consistency over time in sow scores suggests that locomotion scores do not change
352 rapidly over time, so infrequent visits should give a fair representation of the typical locomotion score of
353 a given herd. However, these recommendations are tentative. A larger study specifically addressing the
354 question of sampling methods would be desirable (see e.g. Mullan et al 2009; Main et al 2010).

355 Three aspects of agreement between observers were analysed. Absolute differences between observers
356 will be considered first. Variation of this type is particularly problematic when absolute consistency is
357 required, for example in order to compare the actual level of lameness between different farms or
358 studies. In our study, the farm manager (observer C) used lower scores and more zeroes than other
359 observers and that his scores didn't match or agree as well with those of other observers. While being
360 cautious not to over-interpret this finding from a single observer, it is notable that in dairy cattle,
361 farmers as a group generally underestimate the incidence of lameness in their own herd (Wells et al
362 1993; Whay et al 2003; Rutherford et al 2009; Leach et al 2010). In a study of lameness in sheep, farmers
363 were asked to estimate the prevalence of lameness and then to carry out direct animal-based scoring.
364 Prevalence recorded was higher from animal-based scoring (than from their initial estimate) but scores
365 were still systematically lower than those of a researcher whilst showing good agreement (correlation)
366 with them (King & Green 2011). It is possible that farmers usually categorise low levels of lameness as
367 'normal', because the key thing for them is to identify the level at which the threshold for treatment
368 occurs (King & Green 2011). Alternatively, since farmers primarily spend time with their own animals, on
369 a farm with widespread low levels of lameness, this may become the 'new normal'.

370 A different interpretation of this finding would be that observer C was correct, and that the other less
371 experienced observers over-scored. In particular, the farm manager may have had more experience of
372 the range of 'normal' gaits in older sows (see observers comments above). Farmers and scientists have
373 different perspectives on animal welfare which may affect their interpretation of the same scoring
374 scheme (Hubbard & Scott 2011).

375 Observer agreement in terms of matching of scores will be considered next. The overall kappa statistic
376 suggested fair agreement, although the PABAK statistic was higher suggesting good agreement (Byrt
377 1996). This difference probably arises because kappa does not just reflect 'agreement' but can be
378 affected by bias, where there are systematic differences between observers, and prevalence, where the
379 distribution of scores is not uniform (Byrt et al 1993). Both bias and prevalence were clearly issues in this
380 dataset, and the PABAK statistic is intended to adjust for these. Rutherford et al (2009) reported a mean

381 PABAK of 0.88 (range 0.67 – 0.94) between 3 or more trained observers (after scores were pooled into
382 'sound' and 'lame' categories) in a study of dairy cattle lameness. The PABAK of 0.692 for the six level
383 scoring system in the present study (Table 3) is quite good in comparison. The level of agreement was
384 comparable to a dairy cow locomotion study by Winckler and Willen (2001) in which the proportion of
385 agreement between three observers was 0.68, while the present study found a slightly higher level of
386 0.743 between four observers (Table 3).

387 When analysing agreement at each score level, Kappa and PABAK were lower for score 1, suggesting
388 poorer agreement. Work in sheep (Kaler et al 2009) and cattle (Flower & Weary 2006) has also shown
389 that agreement is better for higher scores distinguishing between the lowest level (normal) and next
390 level up is often the most difficult. Winckler and Willen (2001) found that 62% of disagreements
391 between observers were at this level, and two other studies found much improved agreement after
392 merging the lowest 2 scores and the top 2 (or 3) scores into simpler non-lame and lame categories
393 (Rutherford et al 2009; Brenninkmeyer et al 2007).

394 Finally, I applied methods of agreement based on association or near matching. Again, levels of
395 agreement found were comparable to other studies Flower and Weary (2006) reported an R^2 from
396 regression of 0.69 between 2 observers, which equates to correlation coefficient (r) of 0.83, higher than
397 the best pairwise agreement measured by spearman's ρ (0.781) or Kendall's τ found in this study (0.771;
398 see table 2). Higher levels of association and matching and lower levels of difference between observers
399 than those found here have been reported in a study of sheep lameness (Kaler et al 2009). This may be
400 in part due to the use of video sequences in this study rather than live scoring.

401 Mixed models showed that inter-observer reliability improved over time and then plateaued at around
402 event 4 or 5 ([Figure 3](#)). This suggests that observers showed better agreement with each other with
403 experience, probably because each showed more internal consistency and perhaps also because after
404 each session observers discussed their experience with the method and how to deal with certain
405 examples or borderline cases. Other studies have shown that locomotion scorers show higher levels of
406 agreement with increasing experience (Main et al 2000) or with more training (March et al 2007;
407 Brenninkmeyer et al 2007; Thomsen et al 2008).

408 Overall, methods to analyse agreement in terms of absolute scores, exact matching and association or
409 near matching all showed that observers did not agree perfectly, typically showing 'moderate'
410 agreement. The levels of agreement were largely in line with other studies of this type which rely on
411 visual assessment and ordinal scoring systems. Ways in which such locomotion scoring systems could be
412 improved are discussed further below. It was notable that agreement for score 1 was poorest and
413 observers reported difficulties with the lower end of the scale. Improved kappa values were obtained
414 when data were combined into 0 and 1 vs 2+, suggesting that a simple sound/ lame system would be
415 preferable, as agreement is better (Brenninkmeyer et al 2007; Rutherford et al 2009). The 3 levels in the
416 Welfare Quality protocol for pigs are equivalent to scores 0-2 (normal or altered gait but 'still using all
417 legs'), score 3 ('minimum weight-bearing on affected limb') and scores 4-5 ('no weight-bearing on
418 affected limb or unable to walk', Welfare Quality® 2009). Others have advocated more complex scoring
419 systems to pick out specific types of abnormality (O'Callaghan et al 2003; Flower & Weary 2006). These

420 might be more precise for research purposes or where early diagnosis for intervention is a priority. If
421 simpler systems are more reliable, these might be better for on-farm overall welfare assessment.

422 *What is the ideal locomotion scoring system?*

423 For any locomotion scoring system based on visual scoring by human observers (and indeed for welfare
424 scoring systems in general), the following attributes are desirable: 1) Easy to use and efficient to carry
425 out on a variety of experimental and commercial situations, 2) Objective, unambiguous descriptions of
426 each score (Garner et al 2002), 3) External validation (Knierim & Winckler 2009) for example against foot
427 pathologies (Flower & Weary 2006; KilBride et al 2010; Kaler et al 2011), analgesics (Rushen et al 2007;
428 Flower et al 2008; Danbury et al 2000) or other behavioural measures such as in broiler chickens latency
429 to lie down in shallow water, which is aversive for them, was associated with locomotion scores (Weeks
430 et al 2002). 4) Training before scoring 'in the field', since training and experience increase agreement
431 between observers (March et al 2007; Brenninkmeyer et al 2007). 5) Users (researchers, assurance
432 schemes) should implement ongoing assessment of inter-observer reliability. Another interesting recent
433 idea is to use a modified visual analogue scale, which retains the advantages of ordinal categorical
434 scoring while adding the advantage of capturing some of the variation within categories (Tuytens et al
435 2009).

436 **Conclusions and animal welfare implications**

437 Valid animal-based scoring methods to assess welfare are important both for research and for on-farm
438 assurance purposes. This study suggested that the locomotion scoring system used was promising,
439 although it would benefit from external validation. Sows were moderately consistent over time in their
440 locomotion scores, and older sows had higher scores. In terms of animal welfare, this suggests that
441 lameness is a chronic problem and that measures to treat it are not entirely successful. In terms of
442 animal welfare assessment methodologies, it suggests that infrequent visits to a farm might give a fair
443 picture of the extent of lameness, provided that a representative range of sow parities was sampled,
444 although this requires further validation in a larger study.

445 Inter-observer agreement can be thought of in terms of three complementary approaches: absolute
446 differences, exact matching and association, which are all useful. Of these, the issue of absolute
447 differences is of greatest importance in terms of 'fairness' of welfare assessment (e.g. comparing the
448 prevalence of lameness between individual farms or systems). Inter-observer agreement was moderate
449 and improved with practice, suggesting that training and regular assessment of inter-observer
450 agreement is important to ensure standardisation of data collection methods. Agreement improved
451 when categories were combined: observers found minor locomotor anomalies difficult to classify. As
452 such a simpler system may be preferable for application of welfare assessment on-farm (e.g. Welfare
453 Quality® 2009).

454 **Acknowledgements**

455 SAC is supported by the Research and Science Division of the Scottish Government. Farrowing sow
456 scoring data were collected as part of the Defra-funded project AW0143. Ian Nevison of BioSS gave
457 statistical advice and ran the non-linear mixed models for ordinal data.

458 **References**

- 459 **Anil SS, Anil L and Deen J** 2005 Evaluation of patterns of removal and associations among culling
460 because of lameness and sow productivity traits in swine breeding herds. *Javma-Journal of the American*
461 *Veterinary Medical Association* 226: 956-961
- 462 **Anil SS, Anil L and Deen J** 2009 Effect of lameness on sow longevity. *Javma-Journal of the American*
463 *Veterinary Medical Association* 235: 734-738
- 464 **Baxter EM, Lawrence AB and Edwards SA** 2011 Alternative farrowing systems: Design criteria for
465 farrowing based on the biological needs of sows and piglets. *Animal* 5: 580-600
- 466 **Blokhuis HJ, Jones RB, Geers R, Miele M and Veissier I** 2003 Measuring and monitoring animal welfare:
467 Transparency in the food product quality chain. *Animal Welfare* 12: 445-455
- 468 **Blokhuis HJ, Keeling LJ, Gavinelli A and Serratosa J** 2008 Animal welfare's impact on the food chain.
469 Trends in Food Science & Technology 19, S79-S87.
- 470 **Bokkers EAM, Leruste H, Heutinck LFM, Wolthuis-Fillerup M, van der Werf JTN, Lensink BJ and Van**
471 **Reenen CG** 2009 Inter-observer and test-retest reliability of on-farm behavioural observations in veal
472 calves. *Animal Welfare* 18: 381-390
- 473 **Botreau R, Bonde M, Butterworth A, Perny P, Bracke MBM, Capdeville J and Veissier I** 2007a
474 Aggregation of measures to produce an overall assessment animal welfare. Part 1: a review of existing
475 methods. *Animal* 1: 1179-1187
- 476 **Botreau R, Bracke MBM, Perny R, Butterworth A, Capdeville J, Van Reenen CG and Veissier I** 2007b
477 Aggregation of measures to produce an overall assessment of animal welfare. Part 2: analysis of
478 constraints. *Animal* 1: 1188-1197
- 479 **Botreau R, Veissier I, Butterworth A, Bracke MBM and Keeling LJ** 2007c Definition of criteria for overall
480 assessment of animal welfare. *Animal Welfare* 16: 225-228
- 481 **Botreau R, Veissier I and Perny P** 2009 Overall assessment of animal welfare: strategy adopted in
482 Welfare Quality (R). *Animal Welfare* 18: 363-370
- 483 **Brenninkmeyer C, Dippel S, March S, Brinkmann J, Winckler C and Knierim U** 2007 Reliability of a
484 subjective lameness scoring system for dairy cows. *Animal Welfare* 16: 127-129
- 485 **Byrt T** 1996 How good is that agreement? *Epidemiology* 7: 561
- 486 **Byrt T, Bishop J and Carlin JB** 1993 Bias, Prevalence and Kappa. *Journal of Clinical Epidemiology* 46: 423-
487 429
- 488 **Danbury TC, Weeks CA, Chambers JP, Waterman-Pearson AE and Kestin SC** 2000 Self-selection of the
489 analgesic drug carprofen by lame broiler chickens. *Veterinary Record* 146: 307
- 490 **De Passille AM and Rushen J** 2005 Can we measure human-animal interactions in on-farm animal
491 welfare assessment? Some unresolved issues. *Applied Animal Behaviour Science* 92: 193-209

- 492 **De Rosa G, Tripaldi C, Napolitano F, Saltalamacchia F, Grasso F, Bisegna V and Bordi A** 2003
493 Repeatability of some animal-related variables in dairy cows and buffaloes. *Animal Welfare* 12: 625-629
- 494 **Engblom L, Lundeheim N, Dalin AM and Andersson K** 2007 Sow removal in Swedish commercial herds.
495 *Livestock Science* 106: 76-86
- 496 **FAWC** 2008 Opinion on Policy Instruments of Protecting and Improving Farm Animal Welfare. London,
497 Farm Animal Welfare Council.
- 498 **Flower FC, Sedlbauer M, Carter E, von Keyserlingk MAG, Sanderson DJ and Weary DM** 2008 Analgesics
499 improve the gait of lame dairy cattle. *Journal of Dairy Science* 91: 3010-3014
- 500 **Flower FC and Weary DM** 2006 Effect of hoof pathologies on subjective assessments of dairy cow gait.
501 *Journal of Dairy Science* 89: 139-146
- 502 **Flower FC and Weary DM** 2009 Gait assessment in dairy cattle. *Animal* 3: 87-95
- 503 **Garner JP, Falcone C, Wakenell P, Martin M and Mench JA** 2002 Reliability and validity of a modified
504 gait scoring system and its use in assessing tibial dyschondroplasia in broilers. *British Poultry Science* 43:
505 355-363
- 506 **Gibbens N** 2008 Animal Health 2008 The Report of the Chief Veterinary Officer. London, Department of
507 Food and Rural Affairs.
- 508 **Gilmour AR, Anderson RD and Rae AL** 1987 Variance-Components on An Underlying Scale for Ordered
509 Multiple Threshold Categorical-Data Using A Generalized Linear Mixed Model. *Journal of Animal*
510 *Breeding and Genetics-Zeitschrift fur Tierzuchtung und Zuchtungsbiologie* 104: 149-155
- 511 **Gonzalez LA, Tolkamp BJ, Coffey MP, Ferret A and Kyriazakis I** 2008 Changes in feeding behavior as
512 possible indicators for the automatic monitoring of health disorders in dairy cows. *Journal of Dairy*
513 *Science* 91: 1017-1028
- 514 **Hubbard C and Scott K** 2011 Do farmers and scientists differ in their understanding and assessment of
515 farm animal welfare? *Animal Welfare* 20: 79-87
- 516 **Jensen TB, Bonde MK, Kongsted AG, Toft N and Sorensen JT** 2010 The interrelationships between
517 clinical signs and their effect on involuntary culling among pregnant sows in group-housing systems.
518 *Animal* 4: 1922-1928
- 519 **Kaler J, George TRN and Green LE** 2011 Why are sheep lame? Temporal associations between severity
520 of foot lesions and severity of lameness in 60 sheep. *Animal Welfare* 20: 433-438
- 521 **Kaler J, Wassink GJ and Green LE** 2009 The inter- and intra-observer reliability of a locomotion scoring
522 scale for sheep. *Veterinary Journal* 180: 189-194
- 523 **Keeling L, Forkman B and Veissier I** 2009 Towards a Welfare Quality Assessment System. Welfare
524 Quality.

- 525 **Keen A and Engel B** 1997 Analysis of a mixed model for ordinal data by iterative re-weighted REML.
526 *Statistica Neerlandica* 51: 129-144
- 527 **Kestin SC, Knowles TG, Tinch AE and Gregory NG** 1992 Prevalence of Leg Weakness in Broiler-Chickens
528 and Its Relationship with Genotype. *Veterinary Record* 131: 190-194
- 529 **KilBride AL, Gillman CE and Green LE** 2009 A cross-sectional study of the prevalence of lameness in
530 finishing pigs, gilts and pregnant sows and associations with limb lesions and floor types on commercial
531 farms in England. *Animal Welfare* 18: 215-224
- 532 **KilBride AL, Gillman CE and Green LE** 2010 A cross-sectional study of prevalence and risk factors for foot
533 lesions and abnormal posture in lactating sows on commercial farms in England. *Animal Welfare* 19:
534 473-480
- 535 **King EM and Green LE** 2011 Assessment of farmer recognition and reporting of lameness in adults in 35
536 lowland sheep flocks in England. *Animal Welfare* 20: 321-328
- 537 **Knierim U and Winckler C** 2009 On-farm welfare assessment in cattle: validity, reliability and feasibility
538 issues and future perspectives with special regard to the Welfare Quality (R) approach. *Animal Welfare*
539 18: 451-458
- 540 **Knowles TG, Kestin SC, Haslam SM, Brown SN, Green LE, Butterworth A, Pope SJ, Pfeiffer D and Nicol**
541 **CJ** 2008 Leg Disorders in Broiler Chickens: Prevalence, Risk Factors and Prevention. *PLoS ONE* 3:
- 542 **Leach KA, Whay HR, Maggs CM, Barker ZE, Paul ES, Bell AK and Main DCJ** 2010 Working towards a
543 reduction in cattle lameness: 1. Understanding barriers to lameness control on dairy farms. *Research in*
544 *Veterinary Science* 89: 311-317
- 545 **Lucia T, Dial GD and Marsh WE** 2000 Lifetime reproductive performance in female pigs having distinct
546 reasons for removal. *Livestock Production Science* 63: 213-222
- 547 **Main DCJ, Barker ZE, Leach KA, Bell NJ, Whay HR and Browne WJ** 2010 Sampling strategies for
548 monitoring lameness in dairy cattle. *Journal of Dairy Science* 93: 1970-1978
- 549 **Main DCJ, Clegg J, Spatz A and Green LE** 2000 Repeatability of a lameness scoring system for finishing
550 pigs. *Veterinary Record* 147: 574-576
- 551 **Main DCJ, Kent JP, Wemelsfelder F, Ofner E and Tuytens FAM** 2003 Applications for methods of on-
552 farm welfare assessment. *Animal Welfare* 12: 523-528
- 553 **Main DCJ, Whay HR, Lee C and Webster AJF** 2007 Formal animal-based welfare assessment in UK
554 certification schemes. *Animal Welfare* 16: 233-236
- 555 **March S, Brinkmann J and Winkler C** 2007 Effect of training on the inter-observer reliability of lameness
556 scoring in dairy cattle. *Animal Welfare* 16: 131-133
- 557 **Mullan S, Browne WJ, Edwards SA, Butterworth A, Whay HR and Main DCJ** 2009 The effect of sampling
558 strategy on the estimated prevalence of welfare outcome measures on finishing pig farms. *Applied*
559 *Animal Behaviour Science* 119: 39-48

- 560 **O'Callaghan KA, Cripps PJ, Downham DY and Murray RD** 2003 Subjective and objective assessment of
561 pain and discomfort due to lameness in dairy cattle. *Animal Welfare* 12: 605-610
- 562 **Rushen J and De Passillé AMB** 1992 The scientific assessment of the impact of housing on animal
563 welfare: A critical review. *Canadian Journal of Animal Science* 72: 721-743
- 564 **Rushen J, Pombourcq E and De Passille AM** 2007 Validation of two measures of lameness in dairy cows.
565 *Applied Animal Behaviour Science* 106: 173-177
- 566 **Rutherford KMD, Langford FM, Jack MC, Sherwood L, Lawrence AB and Haskell MJ** 2009 Lameness
567 prevalence and risk factors in organic and non-organic dairy herds in the United Kingdom. *Veterinary*
568 *Journal* 180: 95-105
- 569 **The Council of The European Union** 2001 Council Directive 2001/88/EC of 23rd October 2001 amending
570 Directive 91/630/EEC laying down minimum standards for the protection of pigs.
- 571 **Thomsen PT, Munksgaard L and Togersen FA** 2008 Evaluation of a lameness scoring system for dairy
572 cows. *Journal of Dairy Science* 91: 119-126
- 573 **Tuytens FAM, Sprenger M, van Nuffel A, Maertens W and Van Dongen S** 2009 Reliability of categorical
574 versus continuous scoring of welfare indicators: lameness in cows as a case study. *Animal Welfare* 18:
575 399-405
- 576 **Veissier I, Butterworth A, Bock B and Roe E** 2008 European approaches to ensure good animal welfare.
577 *Applied Animal Behaviour Science* 113: 279-297
- 578 **Weeks CA, Knowles TG, Gordon RG, Kerr AE, Peyton ST and Tilbrook NT** 2002 New method for
579 objectively assessing lameness in broiler chickens. *Veterinary Record* 151: 762-764
- 580 **Welfare Quality®** 2009 *Welfare Quality® Assessment Protocol for Pigs*. Welfare Quality®
581 Consortium: Lelystad, Netherlands
- 582 **Wells SJ, Trent AM, Marsh WE and Robinson RA** 1993 Prevalence and Severity of Lameness in Lactating
583 Dairy-Cows in A Sample of Minnesota and Wisconsin Herds. *Journal of the American Veterinary Medical*
584 *Association* 202: 78-82
- 585 **Why HR, Main DCJ, Green LE and Webster AJF** 2003 Assessment of the welfare of dairy cattle using
586 animal-based measurements: direct observations and investigation of farm records. *Veterinary Record*
587 153: 197-202
- 588 **Winckler C, Brinkmann J and Glatz J** 2007 Long-term consistency of selected animal-related welfare
589 parameters in dairy farms. *Animal Welfare* 16: 197-199
- 590 **Winckler C and Willen S** 2001 The reliability and repeatability of a lameness scoring system for use as an
591 indicator of welfare in dairy cattle. *Acta Agriculturae Scandinavica Section A-Animal Science* 103-107
592

593 **Table 1** The lameness scoring system used in this study. Observers used the integer scores as instructed,
594 on all but 4 occasions when an intermediate score (e.g. 1.5) was recorded.
595

Score	Label	Description
0	Normal	Even strides, rear end sways slightly while walking, pig is able to accelerate and change direction rapidly. Stands normally.
1	Stiff	Abnormal stride length, movements no longer fluent, pig appears stiff, Pig still able to accelerate and change direction. Stands normally.
2	Slight lameness	Shortened stride, lameness detected, Swagger of rear end while walking, no hindrance in pig's agility. Uneven posture while standing.
3	Lame	Pig slow to get up (may dog-sit), Shortened stride, Minimum weight-bearing on affected limb (standing on toes), Swagger of rear end while walking. May still trot and gallop.
4	Limping	Pig reluctant to get up, holds limb off floor while standing, avoids placing affected limb on the floor while moving
5	Downer	Pig unresponsive- does not move and struggles to stand when encouraged to do so.

596

597

Table 2 Inter-observer agreement statistics are shown as follows: *Differences* between observers' scores. Friedman test for effect of the 4 observers on locomotion score, blocked by sow-event. For pairwise comparisons using Sign and Wilcoxon signed-ranks tests: The score for the second observer was subtracted from the score for the first observer (e.g. for A-B difference = A minus B), so +ve difference means the first score was higher (e.g. A's score higher than B's). *Association* between pairs of observers' scores as calculated by Spearman's Rho (ρ) correlation coefficient and Kendall's Tau (τ) correlation coefficient, Kendall's Coefficient of concordance (W) between multiple observers is also given. *Agreement (exact matching)* between observers' scores given by Fleiss's Kappa (κ) for 4 observers and Cohen's Kappa (κ) for pairs of observers. Weighted Kappa is also shown (using a linear decline in weightings further from the diagonal). Finally, Kappa (κ) calculated after re-coding the data to 1-0 form using two different methods is shown. In the first method, all scores of 1-5 were re-coded as 1 ('lame'), and in the second method, scores of 0 and 1 were re-coded as 0 ('sound'), while scores of 2-5 were re-coded as 1 ('lame'). † = $p < 0.1$ * = $p < 0.05$, ** = $p < 0.01$, *** = $p < 0.001$.

	Sign test (-ve : +ve)	Friedman or Wilcoxon test	Spearman's (ρ)	Kendall's W or τ	Kappa (κ)	Weighted Kappa (κ_w)	Kappa (κ) after combining categories (sound = 0, lame = 1+)	Kappa (κ) after combining categories (sound = 0 or 1, lame = 2+)
Overall (A-B-C-D)		102.6***		0.692***	0.443***	0.582	0.532	0.653
A-B	53 : 28**	1199.5*	0.604***	0.591***	0.482***	0.629	0.549	0.725
A-C	8 : 64***	2352.5***	0.536***	0.524***	0.289***	0.457	0.440	0.478
A-D	46 : 28*	1047.5†	0.643***	0.626***	0.498***	0.642	0.599	0.704
B-C	6 : 85***	3863.0***	0.457***	0.448***	0.249***	0.416	0.335	0.574
B-D	24 : 29	754.0	0.781***	0.771***	0.679***	0.777	0.750	0.852
C-D	74 : 5***	160.0***	0.555***	0.541***	0.290***	0.459	0.439	0.488

Table 3 Raw proportions of agreement between observers broken down into overall agreement (all 4 observers agree) and agreement between pairs. Proportions are worked out using n = 498 as the denominator for calculations of overall agreement. For separate score proportions (for 0, 1, 2, 3, 4 & 5), the denominator was worked out as the average number of scores awarded at that level, either overall (scores in **bold**) or by that pair of observers. The Kappa statistic (κ) and the Prevalence Adjusted Bias Adjusted Kappa (PABAK) are also given *in italics* overall and broken down by score.

Agreement		Proportion of Agreement						
	(n / 498)	Overall	0	1	2	3	4	5
<i>Kappa (κ)</i>		<i>0.443</i>	<i>0.532</i>	<i>0.279</i>	<i>0.486</i>	<i>0.458</i>	<i>0.379</i>	<i>1.000</i>
<i>PABAK</i>		<i>0.692</i>	<i>0.834</i>	<i>-0.176</i>	<i>0.040</i>	<i>0.004</i>	<i>-0.20</i>	<i>1.000</i>
Overall (A-B-C-D)	370	0.743	0.862	0.020	0.200	0.170	0.000	1.000
A-B	417	0.837	0.918	0.370	0.667	0.727	0.500	1.000
A-C	426	0.855	0.937	0.097	0.229	0.600	0.000	1.000
A-D	424	0.851	0.932	0.364	0.612	0.429	0.500	1.000
B-C	407	0.817	0.906	0.198	0.333	0.444	0.000	1.000
B-D	445	0.894	0.951	0.641	0.727	0.462	0.667	1.000
C-D	419	0.841	0.929	0.169	0.294	0.167	0.000	1.000

Table 4 Contingency table of scores for observer B and D to illustrate good agreement. Exact matches (on the diagonal, in **bold**) were summed to give the agreement, which was then expressed as a proportion (of 498; see Table 3).

		Observer D Sow locomotion scores						<i>Totals</i>
		0	1	2	3	4	5	
Observer B Sow locomotion scores	0	382	12	1				<i>395</i>
	1	26	41	6	1			<i>74</i>
	2		1	16	3			<i>20</i>
	3			1	3	1		<i>5</i>
	4				1	2		<i>3</i>
	5						1	<i>1</i>
<i>Totals</i>		<i>408</i>	<i>54</i>	<i>24</i>	<i>8</i>	<i>3</i>	<i>1</i>	<i>498</i>

Table 5 Contingency table of scores for observer B and C to illustrate poor agreement. Exact matches (on the diagonal, in **bold**) were summed to give the agreement, which was then expressed as a proportion (of 498; see Table 3).

		Observer C Sow locomotion scores						<i>Totals</i>
		0	1	2	3	4	5	
Observer B Sow locomotion scores	0	390	3	2				<i>395</i>
	1	65	9					<i>74</i>
	2	9	5	5	1			<i>20</i>
	3	1		2	2			<i>5</i>
	4	1		1	1	0		<i>3</i>
	5						1	<i>1</i>
<i>Totals</i>		<i>466</i>	<i>17</i>	<i>10</i>	<i>4</i>	<i>0</i>	<i>1</i>	<i>498</i>

Figure 1 Diagram of dry sow house where the study was carried out, including dimensions (not to scale). Dashed lines indicate barred gates, shown in their normal position. The top left pen indicates the arc of gate swing to show how gates can be temporarily closed to shut sows into the straw bedded area for mucking out and for locomotion scoring. Sows were observed while being moved from one of the sow pens, along the dunging passage, into the service pen and back again.

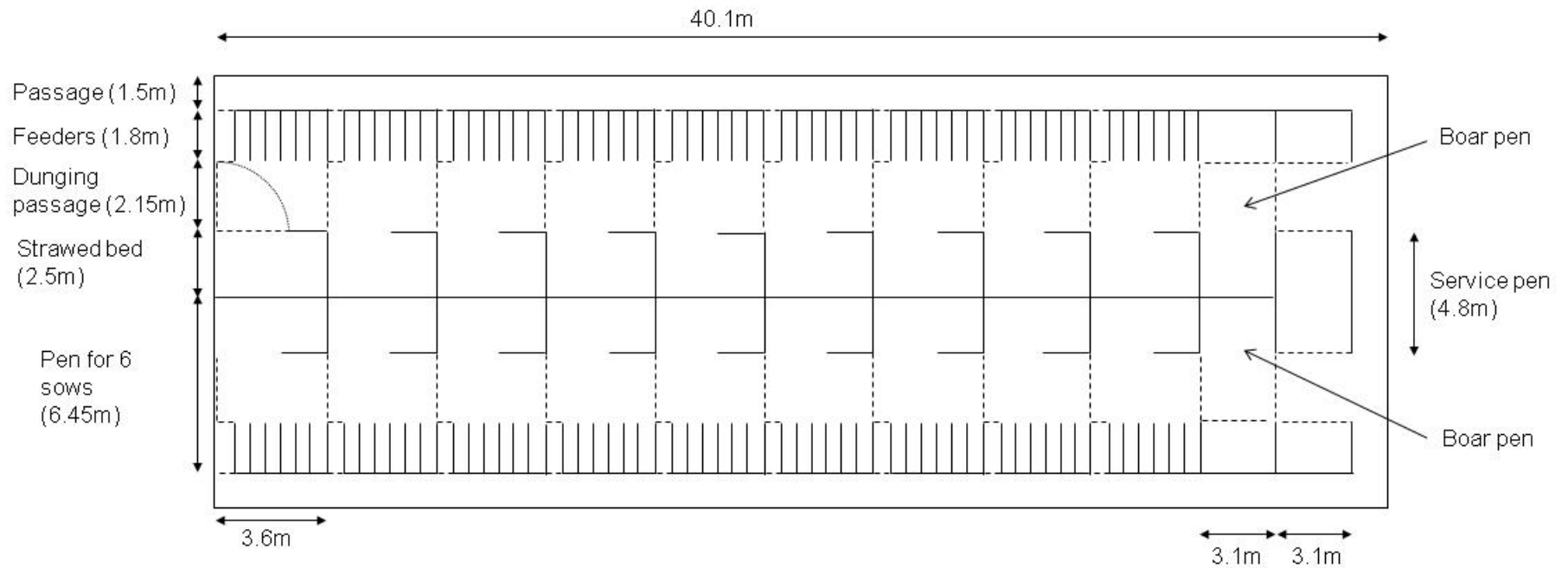


Figure 2 Frequency of locomotion scores given overall by four different observers (A, B, C and D). This data is all scores generated by these 4 observers used in the study (i.e. for all sows at all scoring events), so includes repeat observations of the same sows. Note the broken scale on the Y axis.

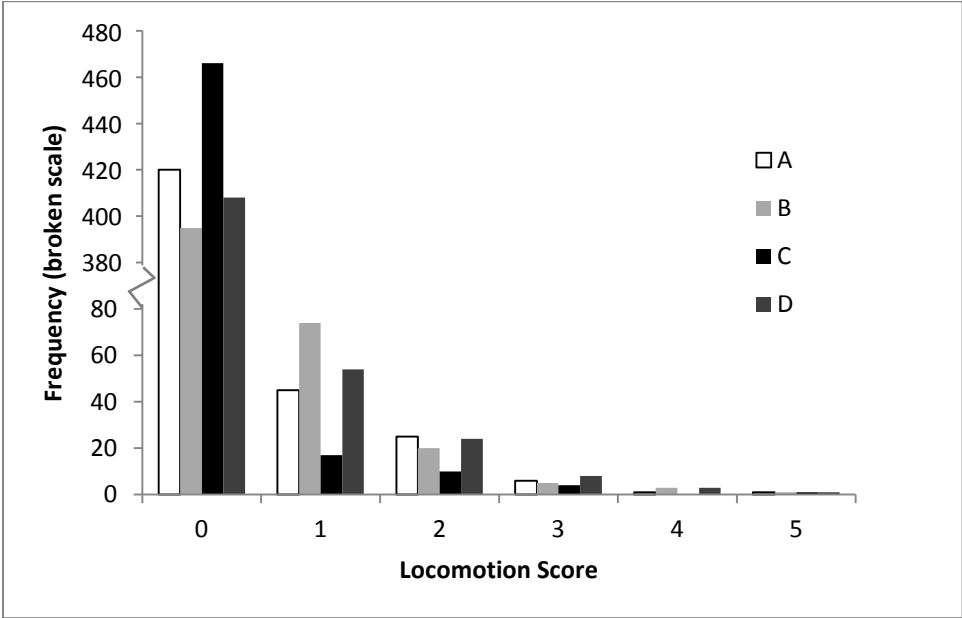


Figure 3 Histogram showing the distribution of sow parities at the end of the study. Sows which have not yet been served are shown as parity 0, parity 1 sows have produced one litter etc. (n=113, data not available for 41 sows).

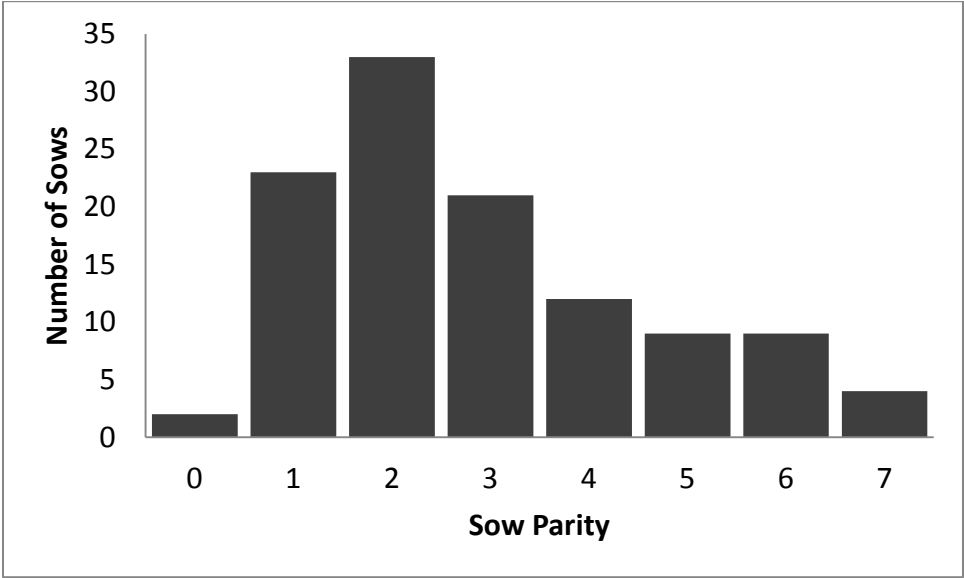


Figure 4 Histogram showing the number of sows at each stage of pregnancy. Sows were usually moved into the farrowing house 109 days after service. N = 563 for this graph as individual sows can appear repeatedly.

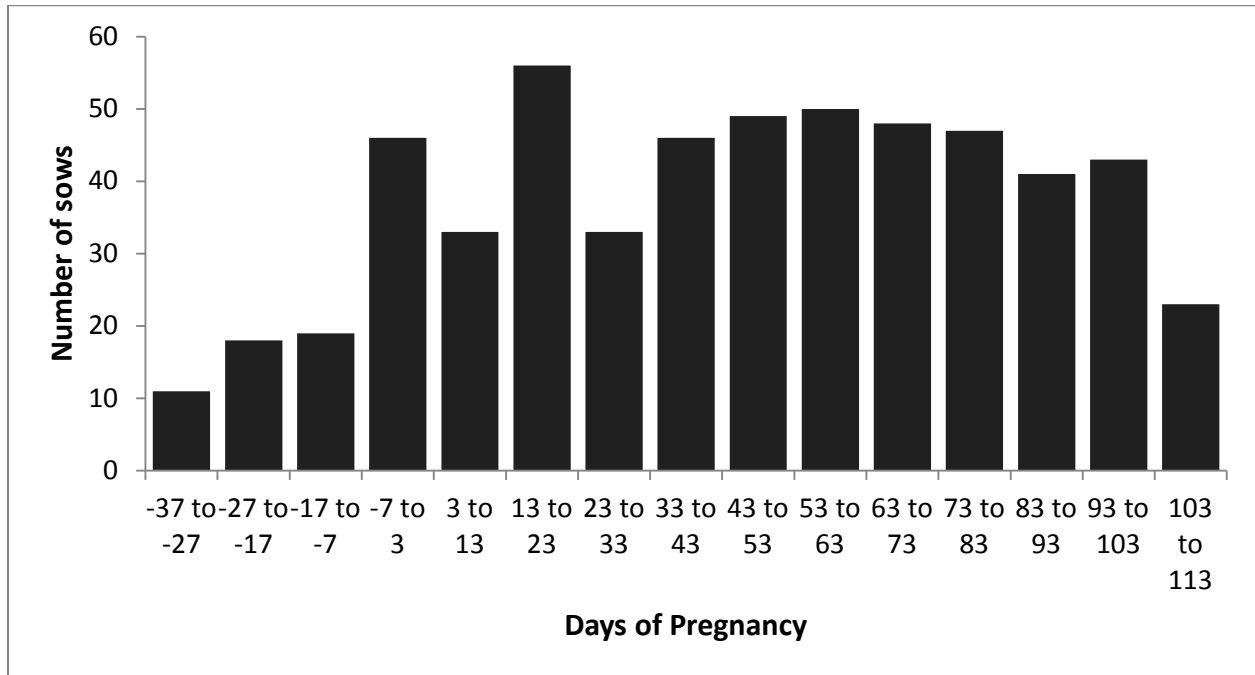


Figure 5 Repeatability of scores across different observers for the 11 scoring events, using data for all observers present at each event. Repeatability estimates were obtained using non-linear mixed models for ordinal data. A separate model was run for each scoring event.

