

Sensitivity of the integrated Welfare Quality[®] scores to changing values of individual dairy cattle welfare measures

S de Graaf^{†‡}, B Ampe[†], S Buijs[†], SN Andreasen[‡], A De Boyer Des Roches^{§‡},
FJCM van Eerdenburg[#], MJ Haskell[¶], MK Kirchner[¶], L Mounier[§], M Radeski[□], C Winckler[□],
J Bijttebier[†], L Lauwers^{†‡}, W Verbeke^z and FAM Tuytens^{*†}

[†] Flanders Research Institute for Agriculture, Fisheries and Food (ILVO), Burgemeester van Gansberghelaan 92, 9820 Merelbeke, Belgium

[‡] Department of Large Animal Sciences, Faculty of Health and Medical Sciences, University of Copenhagen, Groennegaardsvej 8, DK-1870 Frederiksberg, Denmark

[§] Université de Lyon, VetAgro Sup, UMR1213 Herbivores, 69280 Marcy-L'Étoile, France

[#] Department of Farm Animal Health, Faculty of Veterinary Medicine, Utrecht University, 3584 CL Utrecht, The Netherlands

[¶] SRUC, West Mains Road, Edinburgh EH9 3JG, UK

[¶] Department of Veterinary and Animal Sciences, Faculty of Health and Medical Sciences, University of Copenhagen, Groennegaardsvej 8, DK-1870 Frederiksberg, Denmark

[‡] Institut National de la Recherche Agronomique, UMR1213 Herbivores, Equipe Comportement Animal, Robustesse et Approche Intégrée du Bien-Etre, 63122 Saint Genes Champanelle, France

[□] Animal Welfare Center, Faculty of Veterinary Medicine, Ss Cyril and Methodius University in Skopje, Lazar Pop-Trajkov 5-7, 1000 Skopje, Republic of Macedonia

[□] Division of Livestock Sciences, Department of Sustainable Agricultural Systems, University of Natural Resources and Life Sciences, Gregor-Mendel Straße 33, 1180 Vienna, Austria

^z Department of Agricultural Economics, Ghent University, Coupure Links 653, 9000 Ghent, Belgium

* Contact for correspondence and requests for reprints: frank.tuytens@ilvo.vlaanderen.be

Abstract

The Welfare Quality[®] (WQ) protocol for on-farm dairy cattle welfare assessment describes 33 measures and a step-wise method to integrate the outcomes into 12 criteria scores, grouped into four principle scores and into an overall welfare categorisation with four possible levels. The relative contribution of various welfare measures to the integrated scores has been contested. Using a European dataset (491 herds), we investigated: i) variation in sensitivity of integrated outcomes to extremely low and high values of measures, criteria and principles by replacing each actual value with minimum and maximum observed and theoretically possible values; and ii) the reasons for this variation in sensitivity. As intended by the WQ consortium, the sensitivity of integrated scores depends on: i) the observed value of the specific measures/criteria; ii) whether the change was positive/negative; and iii) the relative weight attributed to the measures. Additionally, two unintended factors of considerable influence appear to be side-effects of the complexity of the integration method. Namely: i) the number of measures integrated into criteria and principle scores; and ii) the aggregation method of the measures. Therefore, resource-based measures related to drinkers (which have been criticised with respect to their validity to assess absence of prolonged thirst), have a much larger influence on integrated scores than health-related measures such as 'mortality rate' and 'lameness score'. Hence, the integration method of the WQ protocol for dairy cattle should be revised to ensure that the relative contribution of the various welfare measures to the integrated scores more accurately reflect their relevance for dairy cattle welfare.

Keywords: animal-based welfare indicators, animal welfare, dairy cattle, integrated welfare index, sensitivity analysis, Welfare Quality[®]

Introduction

Accurate welfare assessment is vital for improving animal welfare. In dairy cattle, measures have been developed and validated for a wide variety of both negative and positive aspects of welfare. However, only a few protocols exist that aggregate the scores of multiple welfare measures into one score or index reflecting the overall welfare status of a given herd. Such an overall welfare status score might be used, for example, in the communication with consumers (food-

labelling), as an incentive for on-farm welfare improvements and as a regulative target (Blokhuis *et al* 2010). Examples of schemes that calculate an overall welfare status of dairy cattle are: a protocol by Whay *et al* (2003) based on the 'Five Freedoms' (Farm Animal Welfare Council 1992) that generates a ranking of herds' welfare status; the Animal Needs Index (ANI) which produces an overall welfare score based on integrating mostly resource-based measures (measures of environmental aspects that

Table 1 Principles, the corresponding criteria and measures used in the Welfare Quality® assessment protocol for dairy cattle welfare.

Principles	Criteria	Measures	Aggregation method measures
Good feeding	Absence of prolonged hunger	Body condition score (% very lean animals)	Spline curve fitting
	Absence of prolonged thirst	Availability and cleanliness of water	Decision tree
Good housing	Comfort around resting	Lying down duration; collisions during lying down; on edge/outside of lying area; cleanliness	Converted to ordinal scores, combined in weighted sums and spline curve fitting
	Thermal comfort	No measure for dairy cattle	Decision tree
	Ease of movement	Free stalls or presence of tethering and exercise	
Good health	Absence of injuries	Lameness; integument alterations	Combined in weighted sums, spline curve fitting and Choquet integrations
	Absence of disease	Respiration/digestive diseases; mastitis; mortality; dystocia; downer cows	Converted to ordinal scores, combined in weighted sums and spline curve fitting
	Absence of pain induced by management procedures	Mutilations (dehorning; tail-docking; use of anaesthetics/analgesics)	Decision tree
Appropriate behaviour	Expression of social behaviour	Incidence agonistic interactions	Combined in weighted sums and spline curve fitting
	Expression of other behaviours	Access to pasture	Spline curve fitting
	Good human-animal relationship	Avoidance distance at feeding place	Combined in weighted sums and spline curve fitting
	Positive emotional state	Qualitative Behavioural Assessment	Combined in weighted sums and spline curve fitting

affect welfare) (Bartussek *et al* 2000); and, finally, the Welfare Quality® (WQ) protocol which categorises the overall welfare status of a herd as ‘excellent’, ‘enhanced’, ‘acceptable’ or ‘not classified’ based on a step-wise integration procedure (Welfare Quality® 2009). The current study focuses on the WQ protocol, as this is the only protocol that uses predominantly animal-based measures to calculate an integrated welfare index. Such measures are generally preferred over resource-based measures as the latter tend to reflect risk factors for welfare impairments instead of directly measuring welfare (Blokhuys *et al* 2003, 2010).

In the EU project Welfare Quality® (WQ), protocols for the welfare assessment of the main types of farm animals (cattle, pigs and chickens) were proposed. The dairy cattle protocol describes 33 welfare measures performed on-farm by means of behavioural observations, qualitative behaviour assessment, an avoidance distance test, a management questionnaire, a resource checklist and clinical scoring (Table 1). Subsequently, three steps are used to integrate separate measures into one overall welfare category. Measures are first integrated into criteria scores on a scale of 0–100 which are, in turn, collated into four welfare principles (‘good feeding’, ‘good housing’, ‘good health’ and ‘appropriate behaviour’). These principle scores are then used to determine herds’ overall welfare category (Welfare Quality® 2009). Integration methods are intended to limit compensation of poor scores with better scores on other welfare aspects (Veissier *et al* 2011). Expert opinion of social and animal scientists and stakeholders was used to determine weights for the integration method (Botreau *et al*

2007). Additionally, the protocols were designed with the intention of modifying and updating assessment methods according to advances in animal welfare science (www.welfarequalitynetwork.net/network/45848/7/0/40).

There has been recent discussion about WQ’s measures and integration methods. Some of the measures have been criticised for their poor or undocumented reliability, validity or feasibility (Knierim & Winckler 2009; de Vries *et al* 2013; de Jong *et al* 2015; Tuytens *et al* 2015; de Graaf *et al* 2017). In addition, studies have indicated that a few, resource-based measures have a disproportionately large influence on the overall welfare category (de Vries *et al* 2013; Heath *et al* 2014). Both critical findings may harm the credibility and validity of the WQ protocol in assessing herd welfare. To further examine the functioning of the WQ protocol for dairy cattle, the aim of the current study was to examine: i) if there is variation in sensitivity of integrated outcomes (criteria and principle scores and overall welfare category) to extremely low and high values of measures, criteria and principles; and ii) the reasons for this variation in sensitivity. More specifically, we aimed to critically evaluate whether differences in sensitivity appear to be deliberate and justifiable rather than unintentional side-effects of the complex integration method. To this end, we performed a sensitivity analysis by replacing individual observed values for a given herd with both the theoretically possible and actually observed worst and best values. The latter values were based on a large database of WQ data that reflect a wide range of herd types in Europe, thereby ensuring a substantial but realistic spread in observed values.

Materials and methods

WQ protocol

Only a brief description of the integration method of the WQ protocol for on-farm dairy cattle welfare assessment is given here. The full protocol can be found at <http://www.welfarequalitynetwork.net/>.

Step 1: From measures to criteria scores

Aggregation starts by combining 33 measures into eleven rather than 12 criteria (Table 1), because no data are collected on-farm for the criterion ‘thermal comfort’. Since the recording scales of measures differ, various aggregation methods are used. For categorical measures, decision trees are used resulting in a score between 0–100 where 100 indicates the best possible score. Other measures are converted to ordinal scores where required (eg scores within ‘comfort around resting’ are converted into three categories: normal, moderate problem or serious problem using thresholds [s] for time needed to lie down and percentages of cows for the other measures) and then combined into index values using weighted sums. Spline functions are used to re-weight these sums based on their severity according to expert opinion. Finally, when multiple spline functions were used, Choquet integrals are used to combine these functions into criteria scores on a scale of 0–100 (Botreau *et al* 2007). These algorithmic operators calculate the criteria scores in such a way that a poor score cannot be fully compensated for by a better score in another measure (Botreau *et al* 2007). Consequently, poor scores will have a greater influence on the integrated scores than good scores. Using Choquet integrals, the weight given to each element (measures or criteria) depends on its value relative to the other elements, where the poorest score always gets the highest weight (Botreau *et al* 2008; Welfare Quality® 2009).

Step 2: From criterion scores to principle scores

To integrate criterion scores into principle scores, Choquet integrals are used (Welfare Quality® 2009). The resulting principle scores range from 0 (worst) to 100 (best). Since no data are collected on-farm for the criterion ‘thermal comfort’, this criterion score is replaced with the best score among ‘comfort around resting’ and ‘ease of movement’.

Step 3: From principle scores to overall welfare category

The third and final integration step is from principle scores to overall welfare category. Dairy welfare in a herd is considered ‘excellent’ when it scores > 50 for each principle and > 75 on two of them. When a herd scores > 15 on each principle and > 50 on at least two of them, it is classified as ‘enhanced’. ‘Acceptable’ herds score > 5 for all principles and > 15 for at least three principles. Herds that do not reach the thresholds for the category ‘acceptable’ are considered ‘not classified’ (Botreau *et al* 2009).

Data collection and collation

To reflect the current range present in Europe across various herding systems, pre-existing research datasets of assessments using the WQ protocol for on-farm dairy cattle welfare were collated from seven European research institutes and included data from ten countries. The collected samples were selected by the research institutes to be representative for: i) small-scale dairy herds in Macedonia (n = 12); ii) non-organic and non-tie stall dairy herds in The Netherlands (n = 60) and France (n = 128); iii) random herds with individual somatic cell count data available (SCC, to be able to calculate WQ scores) in Belgium (n = 140), Scotland (n = 16) and Denmark (n = 42); iv) typical herds for the regional low-input herding systems in Romania, Northern Ireland and Spain (n = 30); and v) loose-housed dairy herds with at least 20 cows in Austria (n = 65). The total number of herds in the collated database was 491. To ensure a homogenous integration method for all data, integrated WQ scores were calculated from raw data using a custom-made integration procedure programmed in R 3.2.2 (R Foundation for Statistical Computing, Vienna, Austria). The R integration programme is available on request. The results were checked for coherence with the INRA WQFA webtool (<http://www1.clermont.inra.fr/wq/>), in which WQ measure values can be entered (for dairy cows, fattening pigs, growing pigs and broilers), and WQ criteria, principle and classification scores can be calculated.

Sensitivity analysis

In order to investigate the extent to which values for separate measures affected the criteria and principle scores and the overall welfare category, each herd-level observation for each measure and each herd was replaced one-by-one with both the theoretically possible and the observed (of the entire dataset of 491 herds) worst and best values. This was repeated for individual criteria and principle scores to assess the impact of criteria and principle scores on the overall welfare category. For these calculations, farms that were already in the highest or lowest overall welfare category were excluded. This decision was made because these excluded farms were not able to shift categories, therefore retaining them would give a distorted picture of the results. Subsequently, the median increase and decrease in criteria and principle scores and the percentage of herds that shifted to a lower or higher overall welfare category were quantified for each replacement by the theoretically and observed worst and best values.

For most measures, values that were altered were scored as either percentage of cows (eg % of severely lame cows) or ‘yes’ and ‘no’ (eg for cleanliness of drinkers). However, for some measures (avoidance distance at the feed rack [ADF], lameness and integument alterations) the aggregated measure indexes rather than individual percentages were replaced with worst and best scores. Since these measures together add up to 100% of animals, changing percentages within these could create an impossible situation (ie percentages would add up to over 100%). In addition, the theoretical best score for the

Table 2 Percentages of herds[†] (n = 491) that were downgraded or upgraded one or two overall welfare categories when individual values at measure level (continuous and binary) were replaced with observed worst and best values per measure.

Principles	Criteria, Continuous measures	Observed median, min-max	Observed worst score		Observed best score, % upgraded one category
			% downgraded one category	% downgraded two categories	
Good feeding	<i>Absence of prolonged hunger</i>				
	% of lean cows [‡]	4, 0–88	53	0	5
Good housing	<i>Comfort around resting</i>				
	Mean time needed to lie down (s)	6, 3–20	10	0	6
	% of cows colliding with housing	33, 0–100	5	0	12
	% of cows lying outside of lying area	0, 0–73	11	0	8
	% of cows with dirty flanks	64, 0–100	0	0	7
	% of cows with dirty lower legs	80, 0–100	2	0	7
	% cows with a dirty udder	37, 0–100	2	0	7
Good health	<i>Absence of injuries</i>				
	Lameness index	88, 37–100	6	0	5
	Integument alterations index	53, 0–100	2	0	4
	<i>Absence of diseases</i>				
	Range of all disease measures [‡]	–	1–2	0	0–1
Appropriate behaviour	<i>Expression of social behaviour</i>				
	Head butts per cow per 15 min	0.5, 0–7	13	0	1
	Displacements per cow per 15 min	0.4, 0–5	16	0	4
	<i>Expression of other normal behaviour</i>				
	Number of hours on pasture	7.5, 0–24	9	0	1
	Number of days on pasture	175, 0–365	9	0	1
	<i>Human-animal interaction</i>				
	ADF index	67, 23–100	13	0	6
	<i>Positive emotional state</i>				
	QBA index	0.3, –11–5	24	1	7
	Criteria, Binary measures	% farms with best score			
Good feeding	<i>Absence of prolonged thirst</i>				
	Water flow	82	22	3	3
	Trough length	18	26	1	19
	Number of water bowls		11	1	20
	Drinker cleanliness	76	23	0	8
	At least two drinkers per cow	84	9	0	1
	<i>Ease of movement</i>				
Good housing	Loose or tied housing	93	38	2	3
Good health	<i>Absence of pain induced by management procedures</i>				
	Dehorning method	5	9	0	3
	Tail-docking method	95	8	0	0

[†] Percentages were based on the herds that were actually able to shift one or two categories. For downgrades of one category, n = 482; for downgrades of two categories, n = 174. For upgrades of one category, n = 491.

[‡] As absence of disease contains a very high number of measures with a very small range of shifts, we present only the range here. All separate measures can be found in the Appendix (see supplementary material to papers published in *Animal Welfare*; <https://www.ufaw.org.uk/the-ufaw-journal/supplementary-material>).

Table 4 Percentages of herds[†] (n = 491) that shifted into a different overall welfare category when individual scores were replaced with observed worst and best criteria or principle scores.

Principles, Criteria	Original observed, median, min-max	Observed worst score		Observed best score	
		% farms downgraded one category	% farms downgraded two categories	% farms upgraded one category	% farms upgraded two categories
<i>Good feeding</i>	40, 4–100	64	100	36	1
• Absence of prolonged hunger	70, 3–100	59	0	6	0
• Absence of prolonged thirst	60, 3–100	35	3	30	1
<i>Good housing</i>	54, 6–86	37	0	13	0
• Comfort around resting	27, 0–80	27	0	13	0
• Ease of movement	100, 15–100	27	0	0	0
<i>Good health</i>	34, 8–86	37	0	23	0
• Absence of injuries	35, 4–100	21	0	8	0
• Absence of diseases	40, 12–100	4	0	7	0
• Absence of induced by management procedures	52, 2–100	9	0	3	0
<i>Appropriate behaviour</i>	35, 6–86	37	0	25	0
• Expression of social behaviour	69, 0–100	16	0	5	0
• Expression of other normal behaviour	64, 0–100	9	0	8	0
• Good human-animal relationship	44, 13–100	14	0	8	0
• Positive emotional state	53, 0–93	24	1	7	0

[†] Percentages were based on the herds that were actually able to shift one or two categories. For downgrades of one category, n = 482; for downgrades of two categories, n = 174. For upgrades one category, n = 491; for upgrades of two categories, n = 317.

measures ‘length of drinking trough’ and ‘number of drinking bowls’ depends on the average number of cows in the herd. Therefore, we replaced these with scores that would meet the requirements for all herds in the dataset (10,000 cm for drinking trough length and 100 for number of drinking bowls) as best scores. For the measures of dehorning and tail-docking, we replaced the actual methods used at each herd with the methods which would generate the best (ie no dehorning, no tail-docking, respectively) and the worst score (ie dehorning using surgery with no anaesthetics or analgesics, tail-docking using a rubber band without anaesthetics and analgesics, respectively).

Results

None of the 491 herds were originally (ie before replacement with worst/best scores) in the ‘excellent’ category, 174 (35%) were in the ‘enhanced’ category, 308 (63%) in the ‘acceptable’ category and nine (2%) in the ‘not classified’ category. For eight of the nine, ‘not classified’ herds, classification was due to a ‘good feeding’ principle score below 5 (the threshold for the not-classified category). The median, minimum, and maximum scores are given at the measure (Table 2) and principle and criterion level (Table 4). For several measures, the observed range spanned the entire theoretical range (ie 0–100 for percentages, 0–24 for hours and 0–365 for days). However, for several other measures (18 out of 33), criteria (six out of 12) and principles (three

out of four), the observed data range was narrower than was theoretically possible (Table 2 and Table 3 [see supplementary material to papers published in *Animal Welfare*; <https://www.ufaw.org.uk/the-ufaw-journal/supplementary-material>]). Only 5% of herds were not dehorned or disbudded, 18% were disbudded using caustic paste, 76% using thermocautery, and 1% were dehorned using surgery. Analgesics and/or anaesthetics were used during these procedures in 24 and 60% of the herds, respectively. Only five (*circa* 1%) herds were tail-docked (three by rubber ring and two by surgery). Analgesics were never used during tail-docking whilst anaesthetics were used in two herds.

Sensitivity analysis using observed values: measurement level

Sensitivity of the overall welfare category

When separate measure values were increased to the observed maximum value (ie to the level of the herd that scored best for that specific measure) fewer herds shifted between overall categories than when separate scores were decreased to the observed minimum value (Table 2). For most measures, the highest percentage of shifts between overall welfare categories occurred between the ‘enhanced’ and ‘acceptable’ category (percentage of shifts ranging from 0–99%). However, when some measures (‘% of lean cows’, ‘number of water bowls’, ‘cleanliness of drinker’ and ‘loose

Table 5 Percentages of herds[†] (n = 491) that shifted into a different overall welfare category when scores at the measure, criterion, and principle level[‡] were replaced with theoretically possible worst and best scores.

	Worst score		Best score
	% downgraded one category	% downgraded two categories	% upgraded one category
<i>Measures[†]</i>			
Lameness index [§]	10	0	5
Head butts per cow per 15 min [§]	16	0	1
ADF index [§]	20	0	6
<i>Criteria[†]</i>			
Absence of injuries [#]	29	1	8
Absence of diseases [#]	36	1	7
Absences of pain induced by management procedures [#]	12	0	3
Good-human-animal relationship [#]	23	0	8
<i>Principles[†]</i>			
Good housing [#]	64	100	13
Good health [#]	64	100	23
Appropriate behaviour [#]	64	100	25

[†] Percentages were based on the herds that were actually able to shift one or two categories. For downgrades of one category, n = 482; for downgrades of two categories, n = 174. For upgrades of one category, n = 491.

[‡] Scores from measures, criteria and principles where replacement with theoretical score generated different results than when replaced with observed score.

[§] Theoretical possible worst score was 100, theoretical best score was 0.

[#] Theoretical possible worst score was 0, theoretical best score was 100.

versus tied housing[†]) were increased to the observed maximum level, the highest percentage of shifts to a higher category were between 'not classified' and 'acceptable' (percentage of shifts ranging from 22–100%).

Replacements of measure values only rarely led to negative shifts of more than one category and never to positive shifts of more than one category (Table 2). The effects of replacing a measure often differed greatly, even between measures that belong to the same principle. 'Good health' was the only principle for which changing the values of any of its underlying measures did not result in a substantial (> 10%) effect on herd classification. All measures that were the only measure of a certain criterion caused a relatively high percentage of herds to shift category: '% of lean cows', 'loose or tied housing' and the 'QBA index' when replaced with the worst possible score, with the exception of the 'ADF index'. Although, seemingly combined with many other measures, most measures of the criterion 'absence of prolonged thirst' had a relatively large influence as well. Most upgrades to a higher overall welfare category were achieved by increasing (to the observed maximum levels) 'number of water bowls', 'trough length', and to a lesser extent '% of cows colliding'. Within the two criteria that contained most measures, either sensitivity was very low for all measures ('absence of disease') or sensitivity was greater for those measures that were attributed the highest weight (ie within 'comfort around resting', the measures for resting behaviour are given a higher weight than cleanliness).

Sensitivity of the principles and criteria scores

The sensitivity analysis of the effect of changes in separate measure values on the principle scores and on the criteria scores (Table 3; <https://www.ufaw.org.uk/the-ufaw-journal/supplementary-material>) showed the same pattern as the sensitivity analysis of the overall welfare category. The decrease caused by changing a measure to the lowest observed value was usually greater than the increase caused by changing the same measure to its highest observed value. Exceptions to this trend often concerned measures of which the observed values were very poor. Furthermore, measures that caused the greatest difference tended to belong to criteria that contain few other measures. Exceptions to this trend once again concerned most measures within 'absence of prolonged thirst' and the measure '% of cows colliding with housing'. There was a difference in the sensitivity of the principles and the criteria in that measure values have a more direct influence on criteria scores, and therefore had a greater influence on criteria scores than on principle scores.

Sensitivity analysis using observed values: criteria and principle level

Of all principles, alteration of 'good feeding' led to the highest number of negative as well as positive shifts (Table 4). Moreover, replacing the 'good feeding' score with the lowest observed score in the database caused all 'enhanced' herds to be re-categorised as 'non-classified'.

Table 6 Median (min–max) decrease and increase in principle and criterion scores when measures were replaced with worst and best theoretically possible values.

Principles, Criteria	Measures	Change in principle scores		Change in criteria scores	
		Median decrease in worst scenario	Median increase in best scenario	Median decrease in worst scenario	Median increase in best scenario
<i>Good feeding</i> [†]					
Absence of prolonged hunger	% lean cows [‡]	25 (2–73)	5 (0–69)	69 (2–100)	30 (0–98)
<i>Good health</i> [†]					
Absence of injuries	Lameness index [§]	15 (2–39)	5 (0–35)	27 (3–69)	33 (0–57)
Absence of diseases	Number of coughs per cow per 15 min [‡]	4 (0–12)	0 (0–0)	10 (5–35)	0 (0–0)
	% cows with hampered respiration [‡]	4 (1–12)	0 (0–1)	10 (6–35)	0 (0–14)
<i>Appropriate behaviour</i> [†]					
Good-human-animal relationship	ADF index [‡]	46 (11–82)	9 (0–37)	44 (13–100)	55 (0–87)

[†] Scores from where replacement with theoretical score generated different results than when replaced with observed score.

[‡] Theoretical possible worst score was 100, theoretical best score was 0.

[§] Theoretical possible worst score was 0, theoretical best score was 100.

Alterations to the other principle scores never caused a change of more than one overall welfare category. Alteration of the ‘good housing’ principle caused the fewest positive shifts of all principles, as most farms already scored relatively high for this principle (median score of 54).

Of all criteria, replacement with the lowest observed score was most effective in generating negative shifts for ‘absence of prolonged hunger’ followed by ‘absence of prolonged thirst’. Replacement with the highest observed score was most effective in generating a positive shift for ‘absence of prolonged thirst’. Both criteria within the principle ‘good housing’ (‘comfort around resting’ and ‘ease of movement’) caused 27% of herds to be downgraded when replaced by the observed minimum. Effects of replacing criteria scores within the ‘good health’ and ‘appropriate behaviour’ principles varied considerably between criteria.

Differences between replacement with observed and theoretically possible scores

For several measures, criteria and principles, the observed range did not span the entire theoretical range. For three measures (‘lameness index’, ‘head butts per cow per 15 min’ and ‘ADF index’), four criteria (‘absence of injuries’, ‘absence of diseases’, and ‘absence of pain induced by management procedures’) and three principles (‘good housing’, ‘good health’ and ‘appropriate behaviour’), replacement with the theoretically possible scores instead of the observed scores resulted in a higher % of herds shifting between overall welfare categories (Table 5). For four measures (‘% lean cows’, ‘lameness index’, ‘number of coughs per cow per 15 min’, ‘% cows with hampered respiration’ and ‘ADF index’), this resulted in a higher median increase or decrease of the principle and criteria scores than when worst or best observed scores were used (Table 6).

Discussion

This study investigated the sensitivity of the integrated scores of the WQ protocol for on-farm dairy cattle welfare assessment to extreme changes in individual measure, criterion and principle scores. The impact of one-by-one replacement of observed herd-level measure, criteria and principle scores by extremely low or high values had variable effects on the more highly integrated scores and on the overall welfare category. Investigation into what type of replacements have a large versus a negligible impact suggests that a considerable part of this variation appears to be an unwanted side-effect of the complex step-wise integration method rather than being intentional or justifiable.

Sensitivity analysis using observed values: measurement level

Generally, the impact of a replacement with an extremely low score was bigger than replacement with an extremely high score. This reflects the intention of the WQ integration method to limit compensation of poor scores with better scores on other welfare aspects (Veissier *et al* 2011). The effect of replacing observed measure scores with extreme values on more highly integrated scores (criteria and principles) and on the overall welfare category was very variable and seemed to depend on various aspects. Replacements of the measures ‘% of lean cows’, ‘loose/tied housing’, the ‘QBA index’, ‘drinker trough length’ and ‘cleanliness of drinkers’, had a bigger impact on overall classification compared to other measures (particularly when substituted by observed worst scores). The common feature shared by the first three measures is that they are the only measure of the criterion they belong to (‘absence of prolonged hunger’, ‘ease of movement’ and ‘positive emotional state’, respectively). One other criterion is also documented by a single measure, namely ‘expression of other normal behaviour’

measured with the ADF-test. This measure had less impact compared with the aforementioned three measures, presumably because the ADF-index was already poor for most farms to begin with (so the change by replacing the actual score with the worst possible score was often very small).

The relatively large impact of drinker space and cleanliness of drinkers is in accordance with previous findings for both the dairy cattle protocol (de Vries *et al* 2013; Heath *et al* 2014) and the WQ broiler chicken protocol (Buijs *et al* 2016). This seems to be caused by a combination of factors. First, these measures both belong to the criterion of ‘absence of prolonged thirst’, which contains few measures that matter for calculating the criterion scores (in the decision tree only number/length of drinkers and cleanliness are taken into account). The other measures are either prerequisites for the required number/length of drinkers and therefore less directly influence criterion scores (‘water flow’), or are related to the number of drinkers (‘at least two drinkers per cow’). Second, the principle ‘good feeding’ contains only one other criterion apart from ‘absence of prolonged thirst’, whereas most other principles are composed of more criteria. It could be argued that the large impact of these measures is not necessarily problematic if they are valid indicators of an important welfare problem. However, as resource-based measures, drinker space and cleanliness would appear to be potential risk factors rather than direct measures of thirst (Sprenger *et al* 2009; Vanderhasselt *et al* 2014). Moreover, to our knowledge, the validity of these measures of thirst has not yet been tested. Therefore, the finding that these measures have a relatively large influence on integrated scores can be considered problematic. Animal-based indicators of thirst have been developed, such as blood sodium concentrations, plasma osmolality (Reece 2009; Vanderhasselt *et al* 2013) and voluntary water consumption (in broiler chickens; Sprenger *et al* 2009; Vanderhasselt *et al* 2014). While blood parameters are too invasive to perform in on-farm welfare monitoring, it could be promising to further develop voluntary water consumption tests. Identifying the most reliable, valid and feasible measure of prolonged thirst in dairy cattle should be a priority in future animal welfare assessment research.

Replacements of measures within the principle ‘good health’ with the best or worst scores had little influence on principle and criterion scores and on overall classification, in accordance with previous results (de Vries *et al* 2013; Heath *et al* 2014; Nielsen *et al* 2014). This is remarkable because it includes measures that, according to many experts, indicate important welfare problems in dairy cattle, such as mortality, mastitis and lameness (Lievaart & Noordhuizen 2011; Nielsen *et al* 2014). In addition, Tuytens *et al* (2010) reported that both consumers and farmers rank health aspects as the most important for farm animal welfare. The very limited effect of extreme changes in measures within the criterion ‘absence of diseases’ on integrated WQ scores seems to be caused, at least partially,

by the aggregation method of this criterion. In this aggregation, prevalence of symptoms of diseases is compared to warning and alarm thresholds (eg warning threshold for nasal discharge is 5% of cows and alarm threshold 10% of cows). Subsequently, a weighted sum is calculated of warnings and alarms, with a weight of 1 for warnings and 3 for alarms, which is computed into the criterion score using a spline function. Due to this method, increasing prevalence of diseases that were already above the alarm threshold (or decreasing those that were already below the threshold) will not affect classification at all. Also, when the prevalence of one disease symptom changes, it has only a limited effect on the criterion scores as it is aggregated with many other disease symptoms.

Similarly, measures within ‘absence of injuries’ also had a small impact on the integrated scores. However, a different method is used to integrate the measures within ‘absence of injuries’ to one score. Partial scores for lameness and integument alterations are first calculated using weighted sums and i-spline curves, and are then combined using a Choquet integral. The lameness index had most influence, but still caused only 10% of herds to be downgraded when replaced with the theoretically worst possible score (ie 100% severely lame cows). This surprisingly low impact seems to be due to the large number of criteria within the principle ‘good health’, and to the observation that herds often score relatively low for these criteria. Therefore, changing another score within this principle to a low score is likely to have a smaller effect than when it is done for a score in another principle with fewer criteria such as ‘good feeding’. Due to the limited impact of good health measures on overall welfare categorisation, in theory a situation could occur where farms categorised as ‘acceptable’ or better have 100% severely lame animals, while this may obviously be considered a major welfare problem.

Regarding positive shifts, the percentage of cows colliding with housing had a relatively large positive impact when replaced with best observed score. This is likely because a large proportion of farms (55%) were classified as having a serious problem for this measure to begin with, so for many farms a vast improvement was possible (compared to 37% for ‘% of cows laying out’ and 28% which were above the threshold value of 6.3 s for ‘mean time needed to lie down’).

Sensitivity analysis using observed values: criteria and principle level

There are two, three, or four criteria per principle. This difference in the number of criteria is reflected in the results of the sensitivity analysis: replacement with the worst criteria scores within the principle (‘good feeding’) containing only two criteria (‘absence of prolonged hunger’ and ‘absence of prolonged thirst’) generated most shifts towards a different welfare category. The principle ‘good housing’ also consists of only two criteria for which measures have been developed (for its third criterion ‘thermal comfort’ no measure is available). The impact of both criteria is smaller compared to the two criteria of ‘good feeding’. However, even though for ‘thermal comfort’ no

data are collected, the missing criterion score is replaced with the best score among 'comfort around resting' and 'ease of movement'. This dilutes the effect of a very low score on either of these two criteria. Although some validated measures for thermal comfort exist for dairy cattle (eg respiration rate; Schutz *et al* 2010), inclusion of such measures may complicate timing of farm visits, as the outcomes of these measures are highly influenced by ambient temperature and humidity. Therefore, climatic conditions should be similar during farm visits to capture farm-level differences in thermal comfort rather than differences based on ambient weather conditions. Further research on how to deal with these complexities in the WQ protocol is necessary, or removal of 'thermal comfort' as a criterion for dairy cattle welfare should be considered.

In line with the criteria, of all principles, alteration of 'good feeding' led to the most negative and positive shifts when replaced with observed worst and best scores. For negative shifts this was because 'good feeding' was the only principle for which scores < 5 were observed, which automatically categorises a herd as 'not classified'. For positive shifts, this was because this principle caused more 'not classified' and 'acceptable' categorisations than any other principle (as 131 farms originally had a score between five and 15 for this principle, as opposed to nine for housing, three for health and 23 for behaviour). Therefore, more positive shifts could occur when 'good feeding' was altered than when the other principles were replaced with observed maximum scores.

Differences between replacement with observed and theoretically possible scores

As the sample size in the current study was large and contained a wide variety of herds (given the different sampling aims), we can draw some conclusions about the observed scores in relation to theoretical possible scores. For most measures, observed scores spanned the entire theoretical range. This means that for the dairy cattle protocol, most limits set by WQ seem realistically attainable. For some measures, however, observed scores were less extreme than the theoretically possible scores. In most cases, this did not affect criterion scores as these were within the criterion 'absence of diseases', where warning and alarm thresholds are used to integrate scores. For the lameness index and ADF index, however, fewer shifts of the overall welfare category were observed when replaced with the observed scores. This was also reflected in the corresponding criteria and principle scores, of which the worst possible score never occurred. This is one of the reasons that the principles 'good health' and 'appropriate behaviour' never caused herds to be categorised as 'not classified' when replaced by the observed minimum score.

Animal welfare implications

This study indicates that the WQ integration method does not adequately balance the relative importance of all welfare measures that are included in order to adhere to the multi-dimensional nature of animal welfare. Therefore, using the current integrated WQ scores could lead to a focus on a limited set of (often resource-based) measures which is hard to justify. Since this harms the credibility of the assessment protocol, we recommend a revision of the integration method, so that the relative contribution of the various welfare measures to the integrated scores more correctly reflects their relevance for dairy cattle welfare.

Conclusion

The results of the current study provide insight into the functioning of the integration methods for the dairy cattle WQ protocol. Our findings indicate that the sensitivity of integrated scores to replacement of individual scores by extreme scores is dependent on a number of factors which were intended by the WQ protocol: i) the observed value of the specific measure (or criterion), relative to the values of the other measure in the same criterion (or principle); ii) whether the values were replaced by an extremely low or an extremely high value (more impact of the former); iii) the relative weight WQ attributes to the measures. However, two other factors that were not intended and appear to be unwanted side-effects of the complexity of the step-wise integration method also had considerable influence. These factors were: i) the number of measures that are integrated into criteria and principle scores; and ii) the aggregation method of the measures (eg decision trees or weighted sums). The effect of both integration method and grouping is problematic, as it should be the severity of the welfare problem that affects the overall category. As a result, sensitivity is highest for changes in measures of the 'good feeding' principle, of which a large proportion of the measures are criticised for their validity (ie measures of 'absence of prolonged thirst'). However, measures within the principle 'good health' have the lowest impact although some of these measures are considered to most severely affect dairy cattle welfare. For instance, a farm in the 'acceptable' category or higher could theoretically have 100% severely lame animals. The unwanted side-effects of the current WQ integration methods shown in this study warrant research to develop and evaluate alternative integration methods.

References

- Bartussek H, Leeb C and Held S** 2000 *Animal Needs Index for Cattle*. ANI 35, L/2000. <http://www.bartussek.at/veroeffentlichungen/511134991b0db8204/index.html>
- Blokhuis HJ, Jones RB, Geers R, Miele M and Veissier I** 2003 Measuring and monitoring animal welfare: Transparency in the food product quality chain. *Animal Welfare* 12: 445-455

- Blokhuis HJ, Veissier I, Miele M and Jones B** 2010 The Welfare Quality® project and beyond: Safeguarding herd animal well-being. *Acta Agriculturae Scandinavica* 60: 129-140. <https://doi.org/10.1080/09064702.2010.523480>
- Botreau R, Capdeville J, Perny P and Veissier I** 2008 Multicriteria evaluation of animal welfare at farm level: An application of MCDA methodologies. *Foundation of Computing and Decision Sciences* 33: 1-18
- Botreau R, Veissier I, Butterworth A, Bracke MBM and Keeling LJ** 2007 Definition of criteria for overall assessment of animal welfare. *Animal Welfare* 16: 225-228
- Botreau R, Veissier I and Perny P** 2009 Overall assessment of animal welfare: strategy adopted in Welfare Quality®. *Animal Welfare* 18: 363-370
- Buijs S, Ampe B and Tuytens FAM** 2016 Sensitivity of the Welfare Quality® broiler chicken protocol to differences between intensively reared indoor flocks: which factors explain overall classification? *Animal* 15: 1-10
- de Graaf S, Ampe B and Tuytens FAM** 2017 Assessing dairy cow welfare at the beginning and end of the indoor period using the Welfare Quality® protocol. *Animal Welfare* 26: 213-221. <https://doi.org/10.7120/09627286.26.2.213>
- de Jong IC, Hindle VA, Butterworth A, Engel B, Ferrari P, Gunnink H, Moya TP, Tuytens FAM and Van Reenen CG** 2016 Simplifying the Welfare Quality® assessment protocol for broiler chicken welfare. *Animal* 10: 117-127. <https://doi.org/10.1017/S1751731115001706>
- de Vries M, Bokkers EAM, van Schaik G, Botreau R, Engel B, Dijkstra T and de Boer IJM** 2013 Evaluating results of the Welfare Quality® multi-criteria evaluation model for classification of dairy cattle welfare at the herd level. *Journal of Dairy Science* 96: 6264-6273. <https://doi.org/10.3168/jds.2012-6129>
- Farm Animal Welfare Council** 1992 FAWC updates the Five Freedoms. *Veterinary Record* 17: 357
- Heath CAE, Browne WJ, Mullan S and Main DCJ** 2014 Navigating the iceberg: reducing the number of parameters within the Welfare Quality® assessment protocol for dairy cows. *Animal* 8: 1978-1986. <https://doi.org/10.1017/S1751731114002018>
- Knierim U and Winckler C** 2009 On-farm welfare assessment in cattle: validity, reliability and feasibility issues and future perspectives with special regard to the Welfare Quality® approach. *Animal Welfare* 18: 451-458
- Lievaart JJ and Noordhuizen JPTM** 2011 Ranking experts' preferences regarding measures and methods of assessment of welfare in dairy herds using Adaptive Conjoint Analysis. *Journal of Dairy Science* 94: 3420-3427. <https://doi.org/10.3168/jds.2010-3954>
- Nielsen BH, Angelucci A, Scalvenzi A, Forkman B, Fusi F, Tuytens F, Houe H, Blokhuis H, Sørensen JT, Rothmann J, Matthews L, Mounier L, Bertocchi L, Richard M, Donati M, Nielsen PP, Salini R, de Graaf S, Hild S, Messori S, Nielsen SS, Lorenzi V, Boivin X and Thomsen PT** 2014 Use of animal based measures for the assessment of dairy cow welfare-ANIBAM. *EFSA External Scientific Report*. <https://www.efsa.europa.eu/it/supporting/pub/659e>
- Reece WO** 2009 *Functional Anatomy and Physiology of Domestic Animals*. John Wiley & Sons: Iowa, USA
- Schütz KE, Rogers AR, Poulouin YA, Cox NR and Tucker CB** 2010 The amount of shade influences the behavior and physiology of dairy cattle. *Journal of Dairy Science* 93: 125-133. <https://doi.org/10.3168/jds.2009-2416>
- Sprenger M, Vangestel C and Tuytens FAM** 2009 Measuring thirst in broiler chickens. *Animal Welfare* 18: 553-560
- Tuytens FAM, Federici JF, Vanderhasselt RF, Goethals K, Duchateau L, Sans ECO and Molento CFM** 2015 Assessment of welfare of Brazilian and Belgian broiler flocks using the Welfare Quality® protocol. *Poultry Science* 94: 1758-1766. <https://doi.org/10.3382/ps/pev167>
- Tuytens FAM, Vanhonacker F, Van Poucke E and Verbeke W** 2010 Quantitative verification of the correspondence between the Welfare Quality® operational definition of herd animal welfare and the opinion of Flemish herders, citizens and vegetarians. *Livestock Science* 131: 108-114. <https://doi.org/10.1016/j.livsci.2010.03.008>
- Vanderhasselt RF, Buijs S, Sprenger M, Goethals K, Willemsen H, Duchateau L and Tuytens FAM** 2013 Dehydration indicators for broiler chickens at slaughter. *Poultry Science* 92: 612-619. <https://doi.org/10.3382/ps.2012-02715>
- Vanderhasselt RF, Goethals K, Buijs S, Federici JF, Sans ECO, Molento CFM, Duchateau L and Tuytens FAM** 2014 Performance of an animal-based test of thirst in commercial broiler chicken herds. *Poultry Science* 93: 1327-1336. <https://doi.org/10.3382/ps.2013-03720>
- Veissier I, Jensen KK, Botreau R and Sandøe P** 2011 Highlighting ethical decisions underlying the scoring of animal welfare in the Welfare Quality® scheme. *Animal Welfare* 20: 89-101
- Welfare Quality® Consortium** 2009 *Welfare Quality® Assessment Protocol for Cattle*. Welfare Quality® Consortium: Lelystad, The Netherlands
- Why HR, Main DCJ, Webster AJF and Green LE** 2003 Assessment of the welfare of dairy cattle using animal-based measurements: direct observations and investigation of herd records. *The Veterinary Record* 153: 197-202. <https://doi.org/10.1136/vr.153.7.197>